

# Plot Twist: Uncovering Surprising Event Boundaries in Narratives

Zhilin Wang, Anna Jafarpour, Maarten Sap

University of Washington

{zhilinw, annaja}@uw.edu, msap@cs.washington.edu

## Abstract

When reading stories, people can naturally identify sentences in which a new event starts, i.e., *event boundaries*, using their knowledge of how events typically unfold, but a computational model to detect event boundaries is not yet available. In this study, we characterize and detect sentences with expected or surprising event boundaries in an annotated corpus of short diary-like stories. We train a detection model that combines commonsense knowledge and narrative flow features with a RoBERTa classifier. Our results show that, while commonsense and narrative features can help improve performance overall, detecting event boundaries that are more subjective remains challenging for our model. Upon inspection of the model parameters, we find that sentences marking surprising event boundaries are less likely to be causally related to the preceding sentence, but are more likely to express emotional reactions of story characters, compared to sentences with no event boundary. Additionally, our model performs well in the closely related task of detecting stories with endings that do not follow commonsense from a different corpus. Our results show promise of using models with commonsense knowledge and narrative flow features to shed light on detecting event boundaries in narratives.

## Introduction

When people read stories, they can easily detect the start of new events through changes in circumstances or changes in narrative development, i.e., *event boundaries* (Zacks et al. 2007; Bruni, Baceviciute, and Arief 2014; Foster and Keane 2015; Jafarpour et al. 2019b). These event boundaries can be expected or surprising. For example, in the story in Figure 1 based on crowdsourced annotation, “getting along with a dog who does not generally like new people” marks a *surprising* new event, while “their playing fetch together for a long time” is an *expected* new event.

We aim to study whether machines can detect these surprising or expected event boundaries, using commonsense knowledge and narrative flow features. Characterizing features that are informative in detecting event boundaries can help determine how humans apply expectations on event relationships (Schank and Abelson 1977; Kurby and Zacks 2009; Radvansky et al. 2014; Ūnal, Ji, and Papafragou 2019;

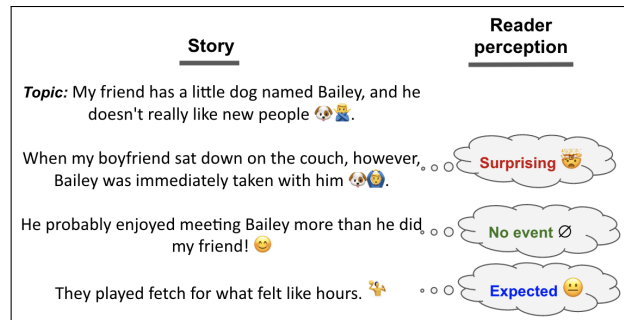


Figure 1: Example story with sentences that contain either a surprising event boundary, no event boundary or an expected event boundary respectively. The annotations of reader perception are from the Hippocampus dataset (Sap et al. 2021).

Zacks 2020). Furthermore, detection of sentences with event boundaries can also be useful when generating engaging stories with a good amount of surprises. (Yao et al. 2019; Rashkin et al. 2020; Ghazarian et al. 2021).

To differentiate sentences with surprising event boundaries, expected event boundaries, and no event boundaries, we train a classifier using 3925 story sentences with human annotation of event boundaries from diary-like stories about people’s everyday lives (Sap et al. 2021). We extract various commonsense and narrative features on relationships between sentences of a story, which can predict the type of event boundaries. Commonsense features include the likelihood that adjacent sentences are linked by commonsense relations from the knowledge graphs Atomic (Sap et al. 2019a) and Glucose (Mostafazadeh et al. 2020). Narrative features include Realis (Sims, Park, and Bamman 2019) that identifies the number of event-related words in a sentence, Sequentiality (Radford et al. 2019; Sap et al. 2021) based on the probability of generating a sentence with varying context and SimGen (Rosset 2020), which measures the similarity between a sentence and the sentence that is most likely to be generated given the previous sentence. We then combine the prediction based on these features with the prediction from a RoBERTa classifier (Liu et al. 2019), to form overall predictions.

We evaluate the performance of the classification model

by measuring F1 of the predictions and compare various configurations of the model to a baseline RoBERTa model. We find that integrating narrative and commonsense features with RoBERTa leads to a significant improvement (+2.2% F1) over a simple RoBERTa classifier. There are also individual differences on the subjective judgment of which sentences contain a surprising or an expected event boundary, that is reflected in the detection model’s performance. The performance of our model increases with increasing agreement across the human annotators. Additionally, by interpreting the trained parameters of our model, we find that the absence of causal links between sentences is a strong predictor of surprising event boundaries.

To further analyze how surprising event boundaries relate to deviation from commonsense understanding, we compare the performance of the classification model on the related task of ROC Story Cloze Test (Mostafazadeh et al. 2016). This task concerns whether the ending sentence of a story follows/violates commonsense based on earlier sentences, which can be linked to whether sentences are expected or surprising. Results show that detecting surprising/expected event boundaries is related to identifying story endings that violate commonsense in the ROC Story Cloze Test, as story endings that violate commonsense are more likely to be predicted to contain surprising event boundaries. However, the performance of our model is significantly higher on the ROC Story Cloze Test (87.9% F1 vs 78.0% F1 on our task). Our findings show that while predicting surprising event boundaries can be correlated with identifying sentences that do not follow commonsense, surprising event boundaries go beyond merely violating commonsense and therefore can be seen as more challenging to detect. Together, our results suggests that while detecting surprising event boundaries remains a challenging task for machines, a promising direction lies in utilizing commonsense knowledge and narrative features to augment language models.

## Event Boundary Detection Task

Events have been widely studied in Natural Language Processing. They have often been represented in highly structured formats with word-specific triggers and arguments (Walker et al. 2006; Li, Ji, and Huang 2013; Chen et al. 2017; Sims, Park, and Bamman 2019; Mostafazadeh et al. 2020; Ahmad, Peng, and Chang 2021) or as Subject-Verb-Object-style (SVO) tuples extracted from syntactic parses (Chambers and Jurafsky 2008; Martin et al. 2018; Rashkin et al. 2018; Sap et al. 2019a). In narratives, events are represented as a continuous flow with multiple boundaries marking new events (Zacks et al. 2007; Graesser, Robertson, and Anderson 1981; Kurby and Zacks 2008; Zacks 2020); however, we lack a model to detect the boundary events that mark the meaningful segmentation of a continuous story into discrete events.

In this work, we study stories from a cognitive angle to detect event boundaries. Such event boundaries relate to our narrative schema understanding (Schank and Abelson 1977; Chambers and Jurafsky 2008; Ryan 2010), commonsense knowledge (Sap et al. 2019a; Mostafazadeh et al. 2020) and world knowledge (Nematzadeh et al. 2018; Bisk et al.

2020). Event boundaries can be surprising or expected based on the knowledge of how a flow of events should unfold. For example, events can be surprising when they deviate from commonsense in terms of what people would predict (e.g., if someone won something, they should not be sad; Sap et al. 2019a). Surprising events can also be low likelihood events (Foster and Keane 2015) such as seeing someone wear shorts outside in winter, or due to a rapid shift in emotional valence between events (Wilson and Gilbert 2008) such as seeing a protagonist being defeated. Importantly, there are individual differences in how humans segment narratives into events (Jafarpour et al. 2019a).

We tackle event boundary detection as a three-way classification task that involves distinguishing surprising but plausible event boundaries in story sentences from expected event boundaries and no event boundaries. To mirror how humans read stories, we predict the event boundary label for a sentence using all of its preceding sentences in the story, as well as the general story topic as context. *Surprising* event boundaries are novel events that are unexpected given their context, such as a dog getting along with someone despite not typically liking new people. *Expected* event boundaries are novel events that are not surprising, such as a person playing a new game with a dog for a long time given that they like each other. In contrast, sentences with *no event* boundary typically continue or elaborate on the preceding event, such as a person liking a dog given that they get along with the dog (Figure 1). Finally, we study the performance of the detection model with respect to the level of agreement in the crowdsourced annotation.

## Event-annotated Data

We use the event-annotated sentences from stories in the Hippocampus dataset to study event boundaries. This dataset contains diary-like stories about real or imagined everyday life experiences, originally released by Sap et al. (2020), of which a subset of 240 stories were annotated at the sentence level by Sap et al. (2021). For the annotation, eight crowdworkers were shown a story sentence by sentence and were asked to mark whether each sentence contained a new surprising or expected event boundary, or no event boundary at all, based on their subjective judgment (Sap et al. 2021). Summarized in Table 1, based on the majoritarian vote, most sentences (57.5%) contain no event boundaries while 16.6% and 13.0% of sentences contains expected and surprising event boundaries, respectively.

Due to the inherent subjectivity of the task, aggregating labels into a majority label yields low agreement (e.g., 61.7% for surprising event boundaries; Table 1). Therefore, at training time, we use the proportion of annotations for each event boundary type as the label instead of the majority vote, because such distributional information is a better reflection of the inherent disagreement among human judgments (Pavlick and Kwiatkowski 2019). At test time, we use the majority vote as a gold label, since measuring performance on distribution modelling is less intuitive to interpret, and subsequently break down performance by agreement level to take disagreements into account.

Majority label	#Samples (%)	% majority agreement (std)
No event	2255 (57.5)	68.1 (13.9)
Expected	650 (16.6)	58.8 (10.6)
Surprising	509 (13.0)	61.7 (11.9)
Tied	511 (13.0)	41.1 (5.7)
Total	3925 (100)	62.2 (15.2)

Table 1: Descriptive Statistics for Event-Annotated sentences. Majority label refers to the most common annotation of a sample from 8 independent annotators. If there is a tie, it is categorized as tied. Majority agreement is the proportion of sample annotations that agrees with the majority label of a sample.

## Event Boundary Detection Model

We first describe informative commonsense and narrative features that we extract for the event boundary detection model. Then, we describe how we integrate these features with a RoBERTa classifier in our model before detailing our experimental setup. Figure 2 depicts an overview of our model.

### Features

We select a collection of commonsense features (Atomic and Glucose relations) and narrative flow features (Realis, Sequentiality and SimGen). A model is trained separately from our main model for Atomic relations, Glucose relations and Realis while models for Sequentiality and SimGen are used without further training. Features of story sentences are extracted as input into the main model. Because language modelling alone might not be sufficient to learn such features (Gordon and Van Durme 2013; Sap et al. 2019a), we provide the extracted features to the model instead of relying on the language models to learn them implicitly.

**Atomic relations** are event relations from a social commonsense knowledge graph containing numerous events that can be related to one another (Sap et al. 2019a). The event relations in this graph consists of:

- Emotional **Reaction**,
- The **Effect** of an event,
- Want** to do after the event,
- What **Needs** to be done before an event,
- The **Intention** to do a certain event,
- What **Attributes** an event expresses.

When an event affects the subject, the feature name is preceded by an x, while if it affects others, it has an o. For example, an xWant of a sentence *PersonX pays PersonY a compliment* is that *PersonX will want to chat with PersonY*, and an oWant is that *PersonY will compliment PersonX back*. We use Atomic relations because surprising event boundaries can involve breaches of commonsense understanding (Bosselut et al. 2019; Sap et al. 2019a;

Mostafazadeh et al. 2020; Gabriel et al. 2021). Furthermore, some Atomic relations (xReact and oReact) concern emotional affect and therefore can be used to capture changes in emotional valence, which can cause events to be seen as surprising (Wilson and Gilbert 2008).

We train an Atomic relation classifier using a RoBERTa-base model (Liu et al. 2019) to classify event-pairs into one of the nine possible relationship labels as well as a None label (to introduce negative samples). We achieved a validation F1 of 77.15%, which is high for a 10-way classification task. We describe training and other experimental details in the Appendix. When making inferences on the event-annotated dataset, we predict the likelihood that a preceding sentence in a story will be related to the current sentence via each of the nine relationship label. Because Atomic relations are directed relations (e.g., *I ate some cake* xEffect *I am full* is different from *I am full* xEffect *I ate some cake*), we also made the reverse inference in case commonsense relations between sentences exist in the reverse direction. Together, 9 forward atomic relation features and 9 reverse features (marked with ‘-r’) are used.

**Glucose relations** are event relations from another commonsense knowledge dataset containing relations between event-pairs in 10 dimensions (Mostafazadeh et al. 2020). Glucose relation features are used to complement Atomic relation features in its coverage of commonsense relations. Dim-1 to 5 are described below while Dim-6 to 10 are the reverse/passive form of Dim-1 to 5 respectively.

- Dim-1: An event that **causes/enables** another event
- Dim-2: An emotion or basic human drive that **motivates** an event
- Dim-3: A **change in location** that enables an event
- Dim-4: A **state of possessing** something that enables an event
- Dim-5: Any other **attribute** that enables an event

Glucose relation classifier was trained on a RoBERTa-base model to classify event-pairs from its annotated dataset into one of ten possible relation labels as well as a None label. We used the *specific* version of Glucose events represented in natural language. As a result, we achieved a validation F1 of 80.94%. Training and other experimental details are in the Appendix. During inference on the Event-annotated dataset, we predict and use as features the likelihood that the current sentence will be related to a preceding sentence via each relation label.

**Realis** events are words that serve as triggers (i.e., head words) for structured event representations (Sims, Park, and Bamman 2019). Realis event words denote concrete events that actually happened, meaning that a higher number of Realis event words suggests greater likelihood of the sentence containing a new event boundary (expected or surprising). We trained a BERT-base model (Devlin et al. 2019) on an annotated corpus of literary novel extracts (Sims, Park, and Bamman 2019), as inspired by Sap et al. (2020). We achieved a validation F1 of 81.85%. Then, we use the trained model to make inference on story sentences in the Event-annotated dataset. Finally, we used the number of Realis

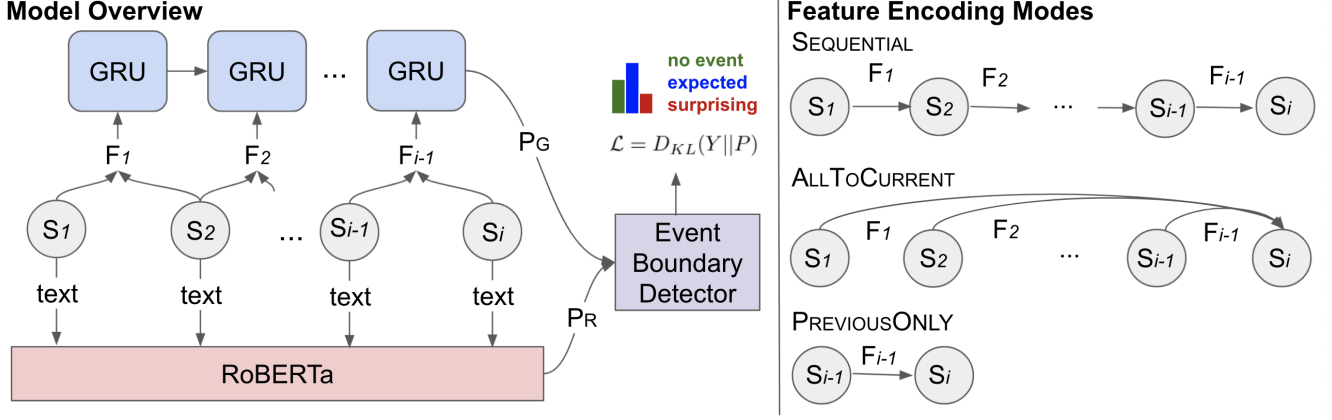


Figure 2: (Left) Our model involves a **GRU** to combine features from sentence pairs with three feature encoding modes, **RoBERTa** to consider story sentences and **Event Boundary Detector** to combine predictions made by the two components.  $S_n$  and  $F_n$  refer to sentence  $n$  and features  $n$  respectively, while  $P_G$  and  $P_R$  are predictions made by the GRU and RoBERTa. The output is a probability distribution over no event boundary, expected event boundary and surprising event boundary, which is used to update model parameters together with the label using the **Kullback-Leibler Divergence** loss function. (Right) **Features** (Atomic, Glucose, Realis, Sequentiality and SimGen) can be extracted as input into the GRU in three feature encoding modes: **SEQUENTIAL** (shown in Model Overview), **ALLTOCURRENT** and **PREVIOUSONLY**.

words in each sentence as a feature. Training and other experimental details are in the Appendix.

**Sequentiality** is a measure of the difference in conditional negative log-likelihood of generating a sentence given the previous sentence or otherwise (Sap et al. 2020, 2021). Sequentiality can be a predictor for unlikely events, which can cause surprise (Foster and Keane 2015). We use GPT-2 (Radford et al. 2019) to measure this negative log-likelihood since it is a Left-to-Right model, which matches the order in which annotators were shown sentences in a story. NLL of each sentence was obtained in two different contexts.  $NLL_{topic}$  is based on the sentence alone with only the topic as prior context, while  $NLL_{topic+prev}$  uses the previous sentence as additional context to study the link between adjacent sentences. Finally, **Sequentiality** is obtained by taking their difference. Experimental details are in the Appendix.

$$NLL_{topic} = -\frac{1}{|s_i|} \log p_{LM}(s_i | Topic)$$

$$NLL_{topic+prev} = -\frac{1}{|s_i|} \log p_{LM}(s_i | Topic, s_{i-1})$$

**SimGen** is computed as the cosine similarity between each sentence and the most likely generated sentence given the previous sentence, under a large Left-to-Right language model (specifically, Turing-NLG; Rosset 2020). Then, we separately converted the original sentence and generated sentence into sentence embeddings using a pre-trained MPnet-base model (Song et al. 2020). Finally, the generated embeddings and the original embeddings are compared for cosine similarity, which is used as a feature. Experimental details are in the Appendix.

## Model Architecture

We propose a model to integrate feature-based prediction with language-based prediction of event boundaries, illustrated in Figure 2 (left). The predictions are independently made with extracted features using a gated recurrent unit (GRU) and with language (i.e., story sentences) using RoBERTa. Then these predictions are combined into a final predicted distribution for the three types of event boundaries. Our model is then trained using the Kullback-Leibler Divergence loss.

**GRU** is used to combine features relating the current sentence  $i$  to prior sentences in a story. It sequentially considers information concerning prior sentences, which mimics the annotator’s procedure of identifying event boundaries as they read one sentence at the time. As seen in Figure 2 (right), we use three feature encoding modes to determine the features that are used as input into the GRU, as inspired by literature on event segmentation (Pettijohn and Radvansky 2016; Baldassano, Hasson, and Norman 2018; Zacks 2020). These three modes represents different ways of facilitating information flow between sentences, which can have distinct effects on identifying event boundaries.

The first mode, **SEQUENTIAL**, encodes features from all previous sentences in the story in a recurrent way (1 to 2, 2 to 3 ...  $i-1$  to  $i$ ) up until the current sentence  $i$ . The second mode, **ALLTOCURRENT**, uses features from each of the previous sentences to the current sentence  $i$  (1 to  $i$ , 2 to  $i$  ...  $i-1$  to  $i$ ). The third mode, **PREVIOUSONLY**, ( $i-1$  to  $i$ ) only feeds into the GRU the features relating to the previous sentence. For all modes, the dimension of each time step input is  $K_G$ , representing the total number of distinct features. We then project the final output of the GRU,  $h_G \in \mathbb{R}^{K_G}$ , into a 3-dimensional vector space representing the unnormalized probability distribution over event boundary types.

**RoBERTa** is used to make predictions based on text in story sentences. We use all story sentences up to sentence  $i$  inclusive. We then project the hidden state of the first token,  $h_R \in \mathbb{R}^{K_R}$ , into a 3-dimensional space representing the unnormalized probability distribution over event boundary types.

**Combining predictions** We combine predictions made by the GRU ( $P_G$ ) and RoBERTa ( $P_R$ ) by concatenating their predictions and multiplying it with a linear classifier of size (6, 3) to output logits of size (3). The logits are then normalized using Softmax to give a distribution of the three types of event boundaries ( $P$ ). The weights of the linear classifier are initialized by concatenating two identity matrix of size 3 ( $\mathbf{I}_3$ ), which serves to perform elementwise addition between the predictions of the GRU and RoBERTa at early stages of the training process.

$$W := [\mathbf{I}_3; \mathbf{I}_3] \quad (1)$$

$$P := \text{Softmax}(W([P_G; P_R])) \quad (2)$$

**Loss function** We use the Kullback-Leibler Divergence loss function to train the model. We use it over the standard Cross Entropy loss function because our training targets are in the form: proportion of annotations for each type of event boundary (e.g., 0.75, 0.125, 0.125 for no event, expected and surprising respectively). Including such distributional information in our training targets over using the majority annotation only can reflect the inherent disagreement among human judgements (Pavlick and Kwiatkowski 2019), which is important to capture for event boundaries given that they are subjective judgements.

## Experimental setup

We seek to predict the event-boundary annotation for each Hippocampus story sentence, using preceding sentences in the story as context, as shown in Figure 2. Additional training and experimental details are available in the Appendix.

**K-fold Cross-validation** Because of the limited size of the dataset ( $n=3925$ ), we split the dataset in k-folds ( $k=10$ ), using one fold ( $n=392$ ) for validation and nine other folds combined for training. From each of the 10 models, we obtained the prediction for the validation set. Together, the validation sets for the 10 models combine to form predictions for the entire dataset, which we use to conduct significance testing in order to compare the performance of models.

**GRU** was accessed from PyTorch, with  $K_G$  set to 33 and a hidden dimension of 33.

**RoBERTa** RoBERTa-base-uncased was used, accessed from HuggingFace Transformers library (with 12-layer, 768-hidden ( $K_R$ ), 12-heads, 110M parameters, 0.1 dropout). When more than 10 prior sentences are available in a story, we use only the most recent 10 sentences due to RoBERTa input sequence length limitations.

	overall F1	no event F1	expected F1	surprising F1
<b>Event Detector (w RoBERTa)</b>				
- PREVIOUSONLY*	78.0	87.2	60.0	59.7
- SEQUENTIAL	77.3	86.6	57.5	60.5
- ALLTOCURRENT	76.9	86.3	57.5	59.7
<b>RoBERTa</b>	75.8	86.2	55.8	54.3
<b>Event Detector (w/o RoBERTa)</b>				
- ALLTOCURRENT	63.9	81.8	32.3	24.8
- SEQUENTIAL	63.8	82.1	34.6	19.5
- PREVIOUSONLY	63.4	81.8	31.8	21.2

Table 2: Event detection task: Performance of Event Detector compared to baseline model. \*: significant difference from RoBERTa based on McNemar’s test ( $p < 0.05$ )

**Evaluation Metrics** While capturing distributional information of subjective judgement labels (Pavlick and Kwiatkowski 2019) is important for training, it can also be difficult to interpret for evaluation. Therefore, we decided to predict for the most likely label during evaluation and compare it against the majority label for each sample. Some samples do not have a single majority label (e.g., equal number of expected and surprising annotations) and these samples were excluded. We use micro-averaged F1 as the metric.

## Results and Discussion

We first quantify the performance of our model in detecting event boundaries, using a coarse-grained performance measure on F1 with respect to majority vote. Then, we investigate how the performance varies based on annotation subjectivity. Finally, we inspect the model parameters to identify commonsense and narrative features that are most informative in detecting event boundaries.

**Improving prediction of event boundaries** As seen in Table 2, RoBERTa alone performs fairly well in predicting event boundaries (F1 = 75.8%, within 2.2% F1 of our best performing model), but can be further supported by our commonsense and narrative features to improve its performance. In contrast, the commonsense and narrative features alone do not perform as well.<sup>1</sup> Overall, our best performing set-up is the Event Detector (PREVIOUSONLY) with F1 = 78.0%, which is significantly different from RoBERTa alone based on McNemar’s test ( $p < 0.05$ ).<sup>2</sup> Its overall strong performance is largely contributed by its strong performance in detecting no event boundaries and expected event boundaries. F1 for no event boundary is higher than for both surprising and expected event boundaries, likely because there are more sentences with no event boundaries as seen in Table 1. The PREVIOUSONLY configuration performs best for no

<sup>1</sup>We also increased learning rate to 1e-3 for better performance given the absence of RoBERTa predictions in this ablation set-up

<sup>2</sup>McNemar’s test is used to determine whether samples that have been predicted accurately (or not) by one model overlap with those that have predicted accurately (or not) by another model

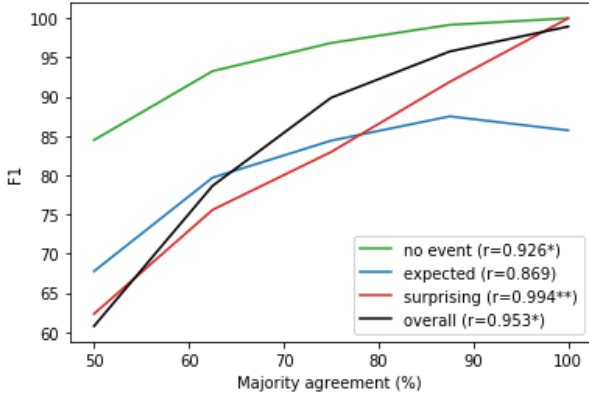


Figure 3: F1 by Event Detector (PREVIOUSONLY) against majority agreement, on all 10 folds. \* means that Pearson’s  $r$  is significant at  $p < 0.05$  and \*\* at  $p < 0.001$ .

event boundaries and expected event boundaries likely because determining whether the current sentence continues an expected event (or not) requires retaining the latest information in working memory (Jafarpour et al. 2019a). However, the SEQUENTIAL configuration seems to perform the best in predicting surprising event boundaries. Compared to no/expected event boundaries, we hypothesize that predicting surprising event boundaries requires taking into account how the story developed prior to the previous sentence in setting up the context for the current sentence. This finding echoes results by Townsend (2018) that showed that surprising sentences take long time to read because it requires changing our mental model formed from previous sentences.

**F1 varies with majority agreement** Since the annotations were subjective and did not always agree, we further examine our best model’s performance (PREVIOUSONLY) with respect to annotation agreement. As shown in Figure 3, F1 increases with majority label agreement (Pearson’s  $r = 0.953$ ,  $p < 0.05$ ). Such positive correlations are observed across all event boundary labels (Pearson’s  $r = 0.869$ – $0.994$ ) and is especially strong for surprising event boundaries (Pearson’s  $r = 0.994$ ,  $p < 0.001$ ). This means that most errors are made on samples that have low agreement among annotators. For example to show this contrast, after “*She and I are very close so it was great to see her marrying someone she loves,*” 7 out of 8 annotators indicated that “*The most memorable moment was when I spilled champagne on my dress before the wedding*” was surprising. On the other hand, after “*It was a hot day in July that our community decided to paint a mural on an intersection for public art,*” only 4 out of 8 annotators indicated that “*I had decided to volunteer to help paint.*” was surprising. The results suggest that our model performance reflects the variability and agreements in humans annotations of event boundaries. We hypothesize that the event boundaries with more agreement are based on features that are shared across the annotators, such as commonsense knowledge; therefore, the model performs well in detecting those. Whereas, our model struggles with detecting event boundaries that are more subjective.

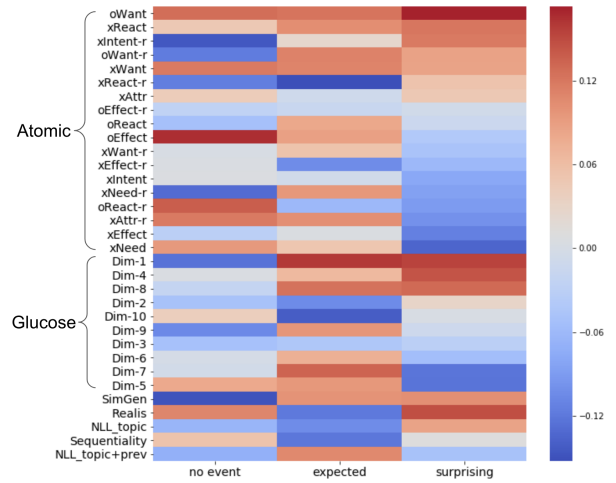


Figure 4: Feature weights towards each label in GRU component of Event Detector (PREVIOUSONLY)

**Predictive features** By integrating a separate feature-based classifier, the Event Boundary Detector model allows us to examine the model parameters and determine features that are associated with surprising, expected or no event boundaries. First, we take the average of the GRU classifier weights for each of the 10 cross-validated models. Then, we plot these weights for each label in Figure 4, and summarize the findings below.

**Features that relate to commonsense relations:** oEffect, xEffect and Glucose Dim-6 (caused by) are most predictive of expected event boundaries. This can indicate that events that are an effect of/caused by a prior event can be expected by annotators, as also noted by Graesser, Robertson, and Anderson (1981). An example of an expected event boundary is “*I told her we could go for coffee sometime.*”, as an effect of “*We had a good time together.*” xNeed is least indicative of surprising event boundaries. This is likely because xNeed refers to what the subject need to do before an activity, which is procedural and unlikely to cause surprise. An example is “*I was grocery shopping a few weeks ago.*” which is needed before “*I had purchased my items and was leaving the store.*”

**Features that explain unlikely events** Realis is highest for surprising event boundaries, suggesting that surprising event boundaries tend to contain the most concrete event-words. Surprising event boundaries also have the highest likelihood when conditioned on the story topic (NLL\_topic) while expected events are highest when conditioned based on the topic and the previous sentence (NLL\_topic+prev). This suggests that surprising events are often inline with the story topic but not with the previous sentence. Therefore, the low likelihood of transitioning between the previous and current sentence is a strong predictor of surprising event boundaries, in line with findings by Foster and Keane (2015) on how the difficulty of linking two adjacent events is an important factor in causing surprise.

**Features that explain changes in emotional valence** Compared to sentences that contain no event boundaries,



sentences that contain either expected or surprising event boundaries have higher `xReact` and `oReact`, which are emotional responses either by the subject or by others to an event. For example, this is the case for the surprising and emotional event boundary *"I remember it was like the 3rd or 4th game when something bad happened."* This suggests that event boundaries are more likely when a sentence is more emotionally charged, echoing work by Dunsmoor et al. (2018) on how event segmentation is particularly frequent when the emotion of fear is triggered.

## Comparison with Story Cloze Test

To better understand how surprising event boundaries relate to deviation from commonsense reasoning, we compare our Event Boundary Detection Task to the ROC Story Cloze Test (Mostafazadeh et al. 2016). This Test involves identifying whether a candidate ending sentence follows commonsense (*commonsense ending*) or deviates from commonsense (*nonsense ending*) given the first four sentences of a short story. Deviation from commonsense reasoning is one factor that can cause surprise (Sap et al. 2019a) and therefore comparing our task to the ROC Story Cloze Test can allow us to potentially isolate deviations from commonsense from other factors that can cause surprise. The ROC Story Cloze Test dataset contains 3142 samples with 1571 commonsense endings and 1571 nonsense endings.<sup>3</sup> For each sample, only one annotation is available. We first train a separate Event Boundary Detector model on the ROC Story Cloze Test. Then, we correlate story endings in the ROC Story Cloze Test with their predicted event boundary label.

	overall F1	nonsense ending F1	commonsense ending F1
<b>Event Detector (w RoBERTa)</b>			
- ALLToCURRENT	87.9	87.8	88.0
- PREVIOUSONLY	87.6	87.3	87.8
- SEQUENTIAL	87.3	87.1	87.5
<b>RoBERTa</b>	87.7	87.6	87.8

Table 3: ROC Story Cloze Test

### Performance of Event Detector on ROC Story Cloze Test

We used the same training and experimental setup as for the event boundary detection task, except the loss function; we use the cross-entropy loss since only one label is available for each sample.<sup>4</sup>

Compared to the Event Boundary Detection task, models perform significantly better on the ROC Story Cloze Test (highest F1 = 78.0% vs. 87.9%,  $p < 0.001$  based on a two-tailed t-test, as observed in Tables 2 and 3). While the tasks are not directly comparable due to the inherent subjectivity

<sup>3</sup>We use the Winter 2018 version, which contains a dev and a test set. As in previous work (Schwartz et al. 2017), we train our model on the dev portion.

<sup>4</sup>Training each model took 20 minutes on an Nvidia P100 GPU.

	no event* (%)	expected (%)	surprising* (%)
nonsense	23.0	19.9	7.1
commonsense	25.3	19.0	5.7

Table 4: Proportion of ROC Story endings that are predicted as each type of event boundary. \*Statistically different between nonsense and commonsense story endings, following a chi-square test with Holm correction ( $p < 0.05$ ).

of the Event Boundary Detection Task, the higher performance on the ROC Story Cloze Test suggests that identifying surprising, expected or no event boundaries may be more challenging than identifying commonsense or nonsense endings. Our commonsense and narrative features also do not seem to significantly improve upon RoBERTa’s performance in the ROC Story Cloze Test (+0.2% F1). This indicates that detecting whether a story ending follows commonsense can be effectively approached using RoBERTa alone, making it relatively easier to tackle.

**Correlation of event boundary with ROC story ending** We use an Event Detector (PREVIOUSONLY) trained on the Event Boundary Detection task to make predictions on the ROC story endings dataset. As seen in Table 4, nonsense endings contain a larger proportion of surprising event boundaries while commonsense endings contains more of no event boundary. This indicates that sentences that breach commonsense are more often surprising, supporting our earlier findings as well as similar conclusion in Graesser, Robertson, and Anderson (1981) and Foster and Keane (2015). On the other hand, sentences that follow commonsense are more often predicted as containing no event boundary, likely as they continue events from prior story sentences. Taken together with the difference in model performance between the two tasks, this suggests that detecting surprising event boundaries is associated with commonsense deviations but is more challenging, likely going beyond merely commonsense breaches.

## Conclusion

We tackle the task of identifying event boundaries in stories. We propose a model that combines predictions made using commonsense and narrative features with a RoBERTa classifier. We found that integrating commonsense and narrative features can significantly improve the prediction of surprising event boundaries through detecting violations to commonsense relations (especially relating to the absence of causality), low likelihood events, and changes in emotional valence. Our model is capable in detecting event boundaries with high annotator agreement but limited in detecting those with lower agreement. Compared to identifying commonsense and nonsense story endings in Story Cloze Test, our task is found to be more challenging. Our results suggest that considering commonsense knowledge and narrative features can be a promising direction towards characterizing and detecting event boundaries in stories.

## References

- Ahmad, W. U.; Peng, N.; and Chang, K.-W. 2021. GATE: Graph Attention Transformer Encoder for Cross-lingual Relation and Event Extraction. arXiv:2010.03009.
- Baldassano, C.; Hasson, U.; and Norman, K. A. 2018. Representation of Real-World Event Schemas during Narrative Perception. *The Journal of Neuroscience*, 38(45): 9689–9699.
- Bisk, Y.; Holtzman, A.; Thomason, J.; Andreas, J.; Bengio, Y.; Chai, J.; Lapata, M.; Lazaridou, A.; May, J.; Nisnevich, A.; Pinto, N.; and Turian, J. 2020. Experience Grounds Language. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 8718–8735. Online: Association for Computational Linguistics.
- Bosselut, A.; Rashkin, H.; Sap, M.; Malaviya, C.; Çelikyilmaz, A.; and Choi, Y. 2019. COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Bruni, L. E.; Baceviciute, S.; and Arief, M. 2014. Narrative Cognition in Interactive Systems: Suspense-Surprise and the P300 ERP Component. In Mitchell, A.; Fernández-Vara, C.; and Thue, D., eds., *Interactive Storytelling*, 164–175. Cham: Springer International Publishing. ISBN 978-3-319-12337-0.
- Chambers, N.; and Jurafsky, D. 2008. Unsupervised Learning of Narrative Event Chains. In *Proceedings of ACL-08: HLT*, 789–797. Columbus, Ohio: Association for Computational Linguistics.
- Chen, Y.; Liu, S.; Zhang, X.; Liu, K.; and Zhao, J. 2017. Automatically Labeled Data Generation for Large Scale Event Extraction. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 409–419. Vancouver, Canada: Association for Computational Linguistics.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186. Minneapolis, Minnesota: Association for Computational Linguistics.
- Dunsmoor, J. E.; Kroes, M. C. W.; Moscatelli, C. M.; Evans, M. D.; Davachi, L.; and Phelps, E. A. 2018. Event segmentation protects emotional memories from competing experiences encoded close in time. *Nature Human Behaviour*, 2(4): 291–299.
- Foster, M. I.; and Keane, M. T. 2015. Predicting Surprise Judgments from Explanation Graphs.
- Gabriel, S.; Bhagavatula, C.; Shwartz, V.; Le Bras, R.; Forbes, M.; and Choi, Y. 2021. Paragraph-level Commonsense Transformers with Recurrent Memory. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(14): 12857–12865.
- Ghazarian, S.; Liu, Z.; S M, A.; Weischedel, R.; Galstyan, A.; and Peng, N. 2021. Plot-guided Adversarial Example Construction for Evaluating Open-domain Story Generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4334–4344. Online: Association for Computational Linguistics.
- Gordon, J.; and Van Durme, B. 2013. Reporting Bias and Knowledge Acquisition. In *Proceedings of the 2013 Workshop on Automated Knowledge Base Construction, AKBC '13*, 25–30. New York, NY, USA: Association for Computing Machinery. ISBN 9781450324113.
- Graesser, A. C.; Robertson, S. P.; and Anderson, P. A. 1981. Incorporating inferences in narrative representations: A study of how and why. *Cognitive Psychology*, 13(1): 1–26.
- Jafarpour, A.; Buffalo, E. A.; Knight, R. T.; and Collins, A. G. 2019a. Event segmentation reveals working memory forgetting rate. Available at SSRN 3614120.
- Jafarpour, A.; Griffin, S.; Lin, J. J.; and Knight, R. T. 2019b. Medial orbitofrontal cortex, dorsolateral prefrontal cortex, and hippocampus differentially represent the event saliency. *Journal of cognitive neuroscience*, 31(6): 874–884.
- Kurby, C.; and Zacks, J. 2009. Segmentation in the perception and memory of events. *Trends in cognitive sciences*.
- Kurby, C. A.; and Zacks, J. M. 2008. Segmentation in the perception and memory of events. *Trends in cognitive sciences*, 12(2): 72–79.
- Li, Q.; Ji, H.; and Huang, L. 2013. Joint Event Extraction via Structured Prediction with Global Features. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 73–82. Sofia, Bulgaria: Association for Computational Linguistics.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv:1907.11692.
- Martin, L.; Ammanabrolu, P.; Wang, X.; Hancock, W.; Singh, S.; Harrison, B.; and Riedl, M. 2018. Event Representations for Automated Story Generation with Deep Neural Nets. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).
- Mostafazadeh, N.; Chambers, N.; He, X.; Parikh, D.; Batra, D.; Vanderwende, L.; Kohli, P.; and Allen, J. 2016. A Corpus and Cloze Evaluation for Deeper Understanding of Commonsense Stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 839–849. San Diego, California: Association for Computational Linguistics.
- Mostafazadeh, N.; Kalyanpur, A.; Moon, L.; Buchanan, D.; Berkowitz, L.; Biran, O.; and Chu-Carroll, J. 2020. GLUCOSE: Generalized and Contextualized Story Explanations. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4569–4586. Online: Association for Computational Linguistics.



- Nematzadeh, A.; Burns, K.; Grant, E.; Gopnik, A.; and Griffiths, T. 2018. Evaluating Theory of Mind in Question Answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2392–2400. Brussels, Belgium: Association for Computational Linguistics.
- Pavlick, E.; and Kwiatkowski, T. 2019. Inherent Disagreements in Human Textual Inferences. *Transactions of the Association for Computational Linguistics*, 7: 677–694.
- Pettijohn, K. A.; and Radvansky, G. A. 2016. Narrative event boundaries, reading times, and expectation. *Memory & Cognition*, 44(7): 1064–1075.
- Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; and Sutskever, I. 2019. Language Models are Unsupervised Multitask Learners.
- Radvansky, G. A.; Tamplin, A. K.; Armendarez, J.; and Thompson, A. N. 2014. Different Kinds of Causality in Event Cognition. *Discourse Processes*, 51(7): 601–618.
- Rashkin, H.; Celikyilmaz, A.; Choi, Y.; and Gao, J. 2020. PlotMachines: Outline-Conditioned Generation with Dynamic Plot State Tracking. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4274–4295. Online: Association for Computational Linguistics.
- Rashkin, H.; Sap, M.; Allaway, E.; Smith, N. A.; and Choi, Y. 2018. Event2Mind: Commonsense Inference on Events, Intents, and Reactions. In *ACL*.
- Reimers, N.; and Gurevych, I. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Rosset, C. 2020. Turing-NLG: A 17-billion-parameter language model by Microsoft.
- Ryan, M.-L. 2010. Narratology and Cognitive Science: A Problematic Relation. *Style*, 44(4): 469–495.
- Sap, M.; Horvitz, E.; Choi, Y.; Smith, N. A.; and Pennebaker, J. W. 2020. Recollection versus Imagination: Exploring Human Memory and Cognition via Neural Language Models. In *ACL*.
- Sap, M.; Jafarpour, A.; Choi, Y.; Smith, N. A.; Pennebaker, J. W.; and Horvitz, E. 2021. Recollection versus Imagination: Exploring Human Memory and Cognition via Neural Language Models. In submission.
- Sap, M.; Le Bras, R.; Allaway, E.; Bhagavatula, C.; Lourie, N.; Rashkin, H.; Roof, B.; Smith, N. A.; and Choi, Y. 2019a. ATOMIC: An Atlas of Machine Commonsense for If-Then Reasoning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01): 3027–3035.
- Sap, M.; Rashkin, H.; Chen, D.; Le Bras, R.; and Choi, Y. 2019b. Social IQa: Commonsense Reasoning about Social Interactions. 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)*, 4463–4473. Hong Kong, China: Association for Computational Linguistics.
- Schank, R.; and Abelson, R. 1977. *Scripts, Plans, Goals, and Understanding*. Hillsdale, NJ: Earlbaum Assoc.
- Schwartz, R.; Sap, M.; Konstas, I.; Zilles, L.; Choi, Y.; and Smith, N. A. 2017. The Effect of Different Writing Tasks on Linguistic Style: A Case Study of the ROC Story Cloze Task. In *CoNLL*.
- Sims, M.; Park, J. H.; and Bamman, D. 2019. Literary Event Detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 3623–3634. Florence, Italy: Association for Computational Linguistics.
- Song, K.; Tan, X.; Qin, T.; Lu, J.; and Liu, T.-Y. 2020. MP-Net: Masked and Permuted Pre-training for Language Understanding. *arXiv preprint arXiv:2004.09297*.
- Townsend, D. J. 2018. Stage salience and situational likelihood in the formation of situation models during sentence comprehension. *Lingua*, 206: 1–20.
- Ünal, E.; Ji, Y.; and Papafragou, A. 2019. From Event Representation to Linguistic Meaning. *Topics in Cognitive Science*, 13(1): 224–242.
- Walker, C.; Strassel, S.; Medero, J.; and Maeda, K. 2006. ACE 2005 Multilingual Training Corpus.
- Wilson, T. D.; and Gilbert, D. T. 2008. Explaining Away: A Model of Affective Adaptation. *Perspectives on Psychological Science*, 3(5): 370–386. PMID: 26158955.
- Yao, L.; Peng, N.; Ralph, W.; Knight, K.; Zhao, D.; and Yan, R. 2019. Plan-And-Write: Towards Better Automatic Storytelling. In *The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)*.
- Zacks, J. M. 2020. Event Perception and Memory. *Annual Review of Psychology*, 71(1): 165–191. PMID: 31905113.
- Zacks, J. M.; Speer, N. K.; Swallow, K. M.; Braver, T. S.; and Reynolds, J. R. 2007. Event perception: a mind-brain perspective. *Psychological bulletin*, 133(2): 273.

## Appendix

### Atomic relations training details

We used the train/dev/test splits from the original Atomic dataset. Negative samples are created by matching a Atomic event node to a corresponding tail event node from another sample based on the relationship involved. Sepcifically, negative sampling was performed on groups ([`'xWant'`, `'oWant'`, `'xNeed'`, `'xIntent'`], [`'xReact'`, `'oReact'`, `'xAttr'`], [`'xEffect'`, `'oEffect'`]) given that the tail event nodes in each group are more similar, creating more discriminating negative samples, as inspired by Sap et al. (2019b). One negative sample is introduced every nine positive samples, since there are nine labels. We used a learning rate of  $1e-4$ , batch size of 64, 8 epochs and AdamW optimizer. Training took 18 hours on a Nvidia P100 GPU.

### Glucose relations training details

Because the Glucose dataset was not split initially, we randomly split the dataset into train/dev/test splits based on a 80/10/10 ratio. For each sample in Glucose, annotations share similar head event nodes in Dim-1 to 5 and similar tail event nodes in Dim-6 to 10. Therefore, our negative sampling strategy for Dim-1 to 5 involves randomly choosing a tail node from Dim-6 to 10 and vice-versa. As a result, one negative sample is introduced every five samples. During training, we used a learning rate of  $1e-4$ , batch size of 64, 8 epochs and AdamW optimizer. Training took 15 hours on a Nvidia P100 GPU.

### Realis training details

We used the train/dev/test split from the Realis dataset. During training, we used the AdamW optimizer, a learning rate of  $2e-5$ , 3 epochs and batch size of 4, as inspired by (Sap et al. 2020). Training took 1 hour on a Nvidia P100 GPU.

### Sequentiality experimental details

GPT2-small was accessed from HuggingFace Transformers library and used without further fine-tuning. It has 125M parameters, a context window of 1024, hidden state dimension of 768, 12 heads and dropout of 0.1.

### SimGen experimental details

We used the Turing-NLG model without further fine-tuning. The model has 17B and we used it with top-p sampling ( $\text{top-p}=0.85$ ), temperature=1.0 and max sequence length of 64 tokens. MPnet-base model was accessed from the Sentence-BERT library (Reimers and Gurevych 2019) and used without further fine-tuning.

### Event Boundary Detection Model training details

AdamW optimizer was used with  $\alpha = 5 * 10^{-6}$ , following a uniform search using F1 as the criterion at intervals of  $\{2.5, 5, 7.5, 10\} * 10^n$ ;  $-6 \leq n \leq -3$ . Learning rate was linearly decayed (8 epochs) with 100 warm-up steps. Batch size of 16 was used. Validation was done every 0.25 epochs during training. Training each model took around 30 minutes on an Nvidia P100 GPU.