

# What is More Likely to Happen Next? Video-and-Language Future Event Prediction

Jie Lei   Licheng Yu   Tamara L. Berg   Mohit Bansal  
Department of Computer Science  
University of North Carolina at Chapel Hill  
{jielei, licheng, tlberg, mbansal}@cs.unc.edu

## Abstract

Given a video with aligned dialogue, people can often infer *what is more likely to happen next*. Making such predictions requires not only a deep understanding of the rich dynamics underlying the video and dialogue, but also a significant amount of commonsense knowledge. In this work, we explore whether AI models are able to learn to make such multimodal commonsense next-event predictions. To support research in this direction, we collect a new dataset, named Video-and-Language Event Prediction (VLEP), with 28,726 future event prediction examples (along with their rationales) from 10,234 diverse TV Show and YouTube Lifestyle Vlog video clips. In order to promote the collection of non-trivial challenging examples, we employ an adversarial human-and-model-in-the-loop data collection procedure. We also present a strong baseline incorporating information from video, dialogue, and commonsense knowledge. Experiments show that each type of information is useful for this challenging task, and that compared to the high human performance on VLEP, our model provides a good starting point but leaves large room for future work.<sup>1</sup>

## 1 Introduction

Given a video clip (*premise event*), humans can often describe logical events that might happen next (*future events*), and interestingly people tend to agree on which future events are more likely to happen than others. Making such predictions requires not only a deep understanding of the rich dynamics underlying the video and dialogue, but also a significant amount of *multimodal commonsense* knowledge about the world. In Figure 1 (*top*),

<sup>1</sup>Dataset, code are available at <https://github.com/jayleicn/VideoLanguageFuturePred>



Figure 1: Video event prediction examples. Given a video (with dialogue) and two future events, the task is to predict which event is more likely to happen following the video. *Top*: an example with a TV show clip. *Bottom*: an example with a YouTube Lifestyle Vlog clip. The correct answer is shown in bold and green. *Premise Summary* and *Rationale* are included for illustration purpose only, they are hidden for the task.

we show an example where commonsense knowledge about inter-human relationships is required, i.e., that a detective typically does not hand evidence to a suspect in a criminal case. Given this knowledge, it is more likely that Beckett (the detective) will take the phone (the evidence) and read the text, than hand the phone to Dean (the suspect).

In this work, we propose *Video-and-Language Event Prediction* (VLEP), a novel dataset and task for fine-grained future event prediction from videos. Given a video with aligned dialogue, and two possible future events, the AI system is required to understand both visual and language semantics from this

video, and commonsense world knowledge, and then make a sound and practical judgment about the future, by choosing the more likely event from two provided possible future events. The VLEP dataset contains 28,726 examples from 10,234 short video clips. Each example (see Figure 1) consists of a *Premise Event* (a short video clip with dialogue), a *Premise Summary* (a text summary of the premise event), and two potential natural language *Future Events* (along with *Rationales*) written by people. These clips are on average 6.1 seconds long and are harvested from diverse event-rich sources, i.e., TV show and YouTube Lifestyle Vlog videos.

Collecting such a dataset is a non-trivial task, as crowd-workers may write trivial negatives (less-likely events) that contain *biases* or *annotation artifacts* (Gururangan et al., 2018), such as negation (e.g., ‘*says nothing*’) or impolite actions (e.g., ‘*hit someone in the face*’), as shown in Table 1. To mitigate this, we combine two recent effective approaches, adversarial human-and-model-in-the-loop data collection (Nie et al., 2020) and adversarial matching (Zellers et al., 2019a), to build a larger, more-challenging, and less-biased dataset. Specifically, 50% of the examples in VLEP are directly annotated by humans over two rounds: round one of standard data collection, i.e., crowd-workers perform the annotations with no model feedback, and round two of adversarial data collection, i.e., crowd-workers perform the annotations with the goal of fooling our basic models trained on round one data (thus avoiding obvious biases). Our analysis shows that the adversarial data collection helps to mitigate dataset bias (reduce trivial negatives), i.e., we notice that a premise-oblivious model (that does not see the premise event) performs worse on data collected on round two than that of round one. Another 50% of the examples are obtained by performing adversarial matching on the human-annotated positive events (more-likely events), i.e., for each premise event, we sample a positive from other premises as a negative, such that the sampled negative is relevant to the current premise while not being overly similar to the true positive. Overall, our dataset is collected via 3 methods (standard-human, adversarial-human, adversarial-matching), hence maintaining a balance between easy and hard examples while reducing potential biases.

To provide a strong baseline for this challenging multimodal future-prediction task, we propose a transformer-based model to incorporate both visual

and textual information from the premise event. We also inject commonsense reasoning knowledge into our model from the ATOMIC dataset (Sap et al., 2019). Our ablation study shows that each part of our model, i.e., video understanding, dialogue understanding, and commonsense knowledge, is useful for the multimodal event prediction. Though our model has shown promising results, it is still not comparable to human performance (67.46% vs. 90.50%), indicating the challenging nature of the multimodal event prediction task and the large scope for interesting future work on our VLEP dataset.

To summarize, our contributions are 3-fold: (1) We propose a new task, Video-and-Language Event Prediction, which requires a model to make fine-grained, multimodal prediction regarding which future event is more likely to happen following a premise video. (2) We introduce a new dataset VLEP for the task, and use two approaches to gather natural hard-negative future-events: adversarial data collection and adversarial matching. This helps mitigate potential annotation artifacts and biases in the dataset. A detailed analysis of VLEP is provided. (3) We present a strong baseline method to benchmark the proposed dataset, and show that incorporating commonsense knowledge improves performance, indicating future directions for this new task (with a large model-human performance gap).

## 2 Related Work

**Video-and-Language Understanding.** Various datasets and tasks have been introduced in this area, such as video captioning (Xu et al., 2016; Rohrbach et al., 2017; Wang et al., 2019; Lei et al., 2020c), video QA (Tapaswi et al., 2016; Jang et al., 2017; Lei et al., 2018), and moment retrieval (Hendricks et al., 2017; Gao et al., 2017a; Lei et al., 2020c). Recently, Liu et al. (2020) propose the video-and-language inference task where a model needs to infer whether a statement is entailed or contradicted by a video. While this task requires judging a statement’s verification w.r.t. existing events, our task requires predicting future events.

**Commonsense Reasoning.** Recently, commonsense reasoning has emerged as an important topic in both language (Zellers et al., 2018, 2019b; Sap et al., 2019) and vision (Vedantam et al., 2015b; Zellers et al., 2019a; Zadeh et al., 2019; Fang et al., 2020) communities. Zellers et al. (2018, 2019b)

build multiple-choice QA datasets for commonsense inference with text context, Zellers et al. (2019a); Park et al. (2020) propose datasets for commonsense-based QA and captioning on still images, Fang et al. (2020) augment MSRVT (Xu et al., 2016) videos with commonsense captions and QAs. In this work, we focus on a more complex type of context (video with dialogue) and a future prediction task, posing challenges for both video-and-dialogue understanding and commonsense reasoning .

**Video Forecasting.** Predicting the future is one of the popular research areas in the vision community. It covers a wide spectrum of topics, including predicting future frames (Vondrick et al., 2016b; Liang et al., 2017), future action labels (Vondrick et al., 2016a; Gao et al., 2017b; Shi et al., 2018; Epstein et al., 2020), future human motions (Fragkiadaki et al., 2015; Mao et al., 2019), etc. While these works mostly study low-level vision or semantic concepts prediction (e.g., action labels), we focus on predicting high-level future events from video and dialogue.

**Bias in Datasets.** It is known that *biases* or *annotation artifacts* (Goyal et al., 2017; Gururangan et al., 2018; McCoy et al., 2019; Tsuchiya, 2018; Poliak et al., 2018; Zellers et al., 2019a) exist in standard human annotated datasets (Bowman et al., 2015; Williams et al., 2018; Antol et al., 2015; Tapaswi et al., 2016; Jang et al., 2017; Kim et al., 2017; Lei et al., 2018). For example, negation words such as *nobody*, *no* and *never* are strong indicators of contradictions (Gururangan et al., 2018) in MNL (Williams et al., 2018). Such superficial patterns are easy for models to exploit, resulting in an overestimate of task performance (Goyal et al., 2017; Gururangan et al., 2018). Zellers et al. (2019a) propose *Adversarial Matching* to mitigate biases in QA, where positive answers are recycled to serve as negatives for other questions. Nie et al. (2020) propose a *Human-And-Model-in-the-Loop Entailment Training* (HAMLET) adversarial data collection strategy to gather challenging examples for NLI. In this work, we adopt both approaches to construct a less-biased and more challenging dataset for the multimodal video+dialogue setting.

### 3 Dataset

The VLEP dataset contains 28,726 examples from 10,234 TV show and YouTube Lifestyle Vlog video

clips. Of these, 50% are collected directly from human annotators over two rounds: (1) round one: standard data collection; (2) round two: adversarial data collection. We collect human examples using Amazon Mechanical Turk (AMT), with an average cost of \$1.10 per example. More detail about the annotators and quality checks are presented in Appendix Section A.2. The other 50% are obtained from human-annotated examples via Adversarial Matching (Zellers et al., 2019a). Hence, overall we build our dataset with 3 collection methods (standard-human, adversarial-human, adversarial-matching), allowing a balance between easy and hard examples while reducing potential biases.

#### 3.1 Video and Language Source

VLEP is built using videos (with English dialogues) from two sources: TV shows and YouTube Vlogs. Both types of videos contain rich physical interactions and dialogues between people and are thus ideal sources for collecting interesting events. We do not use videos from sources like ActivityNet (Caba Heilbron et al., 2015) since they do not have dialogues and typically contain fewer events.

**TV Show Videos.** We use TV show clips from TVQA (Lei et al., 2018). The clips are typically 60-90 seconds long, and are from 6 popular TV shows of 3 genres: 1) sitcom: *The Big Bang Theory*, *How I Met Your Mother*, *Friends*, 2) medical drama: *Grey’s Anatomy*, *House*, 3) crime drama: *Castle*.

**YouTube Lifestyle Vlogs.** While TV shows are good video sources with rich inter-human interactions, they may focus more on scripted content (Lei et al., 2020b). Thus, we also collect a set of YouTube lifestyle vlogs as additional sources, which are typically more natural and live interactive. We first manually identify a list of YouTube channels that contain videos with rich human interactions and dialogues (in English). We filtered out those channels with instructional videos (Miech et al., 2019) or routine videos (Ignat et al., 2019; Fouhey et al., 2018), as they focus more on a single person performing actions, while we desire videos with richer multi-person interactions and dialogues. In addition, the actors in instructional or routine videos typically follow rigid steps (e.g., in cooking videos, they usually follow recipes) to finish a particular task, making it much easier to predict the future events. In the end, we identified 9 channels that contain a diverse set of lifestyle vlog videos on various topics: *travel*, *food*, *daily life*

and *family*, etc. We downloaded all videos from these channels that are published after 2017, which were then verified to ensure high quality. The resulting pool contains 971 videos of 10-30 minutes long. Each video is associated with aligned dialogue text, either written by humans or generated from YouTube’s Automatic Speech Recognition (ASR) system. We segment the videos into 60-second clips. For each video, we drop the first and the last clip, as most of them are high-level introductions or subscription pleas.

### 3.2 Round One: Standard Data Collection

As our task is video event prediction, we aim to collect a set of videos annotated with future event pairs (i.e., *more-likely events* and *less-likely events*, also referred to as *positive events* and *negative events*) that are likely to happen right after the ‘premise’ video. Each event is written in natural language (English), and we require the positive event to be more likely to happen than the negative event.

With this goal in mind, we create our first annotation task on AMT. We present workers (human *writers*) with a 60-90 seconds long video with aligned dialogue subtitle, to encourage them to write events that are related to both the visual content and the dialogue. Workers are required to first select an interesting event from the video with timestamps, similar to previous works (Lei et al., 2018, 2020c). This event is defined as the *premise event*. We also require workers to write a *premise summary* – a natural language (English) description summarizing the premise event. Following Lei et al. (2018, 2020c), for referring a specific person in the video, workers are instructed to either use the character names (e.g., ‘Sheldon’) if they are available in the dialogues or provide a referring expression (Kazemzadeh et al., 2014) (e.g., ‘the man in blue top’) that uniquely refers to a person in the video. Next, given the premise event, workers are required to write two *future events*, one more likely (>50% chance) to happen after the premise event, and one less likely (<50% chance). For example, in Figure 1, the correct answers are the more-likely while the wrong answers are the less-likely. To encourage workers to write more reasonable future event that ground to the premise event,<sup>2</sup> we also require them to provide a *rationale* as to why it is more or less likely. As it is not the

<sup>2</sup>Otherwise, workers sometimes write random events that are not related to the given premise.

---

<b>Type:</b> Negation
<b>Premise Summary:</b> Amy picks up her phone and reads a text message.
<b>More-likely:</b> Amy tells her friends what the text message says.
<b>Less-likely:</b> Amy says <b>nothing</b> at all to her friends.

---

<b>Type:</b> Impolite Actions
<b>Premise Summary:</b> Chandler finds out that Joey used his toothbrush.
<b>More-likely:</b> Chandler starts arguing with Joey for using his toothbrush.
<b>Less-likely:</b> Chandler <b>hits Joey in the face with a punch</b> .

---

Table 1: Example annotation artifacts in the negative future events (*Less-likely events*).

focus of this work, we will release these rationales to support research on textual explanation generation/classification tasks (Huk Park et al., 2018; Zellers et al., 2019a).

Each collected example is verified by three human *verifiers*, by ranking the future events conditioned on the premise event. We only accept an example if at least three out of four users (one writer + three verifiers) reach an agreement, as Hendricks et al. (2017); Nie et al. (2020). In addition, we also discard examples if one of the verifiers thinks the events are against our instructions (e.g., wrong person reference). In total, we collected 6,458 verified examples from 2329 TV show clips. We split them into 70% training, 15% development, and 15% testing splits such that the videos and their corresponding examples only appear in one split.

### 3.3 Round Two: Adversarial Data Collection

While being efficient in data collection, we found the collected negative events in round one are sometimes simple and contain *biases* or *annotation artifacts* (Gururangan et al., 2018). In Table 1, we show typical examples of annotation artifacts. For example, we found workers tend to use negation when writing the less-likely event. This particular type is similar to the *visual priming bias* (Zhang et al., 2016) for *yes/no* questions in VQA (Antol et al., 2015) and the *negation word bias* (Gururangan et al., 2018) in MNLI (Williams et al., 2018). To quantitatively study the effect of these annotation artifacts, we fine-tune a RoBERTa-base (Liu et al., 2019) model to classify which event is more likely to happen, with only the future events from round one’s training data, i.e., the model has no access to the premise event. On round one’s Dev. split, this premise-oblivious model obtains 75.34% accuracy, which is much higher than chance (50%).

Hence, in order to collect harder and less-biased negatives, we make use of an adversarial collection procedure (see Figure 2), in a human-and-model-in-the-loop process (Nie et al., 2020), where models

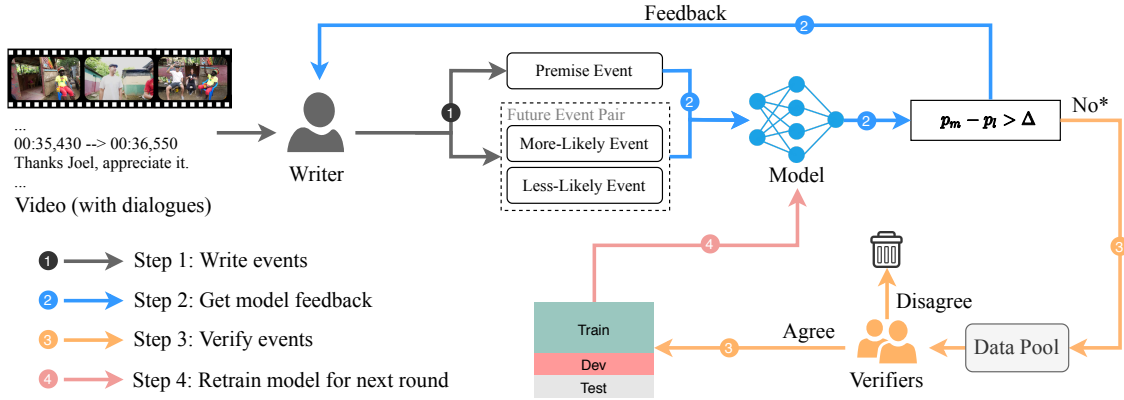


Figure 2: Illustration of our adversarial data collection procedure.  $p_m$  and  $p_l$  are the probabilities of the more-likely and the less-likely event being happening, respectively.  $\Delta$  is a hyperparameter that controls how hard we want the collected example to be, it also helps to reduce prediction noises from imperfect models.  $\Delta$  is set to 0.1 in our experiment. No\* or the number of trials reaches the maximum limit of three.

are used to provide real-time feedback to crowdworkers during data collection. Specifically, each submitted result is sent to the model for evaluation and writers are prompted<sup>3</sup> to rewrite their negative event if our model predicts a much higher probability for the more-likely event ( $p_m$ ) than the less-likely event ( $p_l$ ), i.e.,  $p_m - p_l > \Delta$ , where  $\Delta$  is a hyperparameter that controls how difficult we want the collected examples to be and is set to 0.1 empirically. This can be seen as a *soft-adversarial* strategy, unlike Nie et al. (2020) where feedback decisions are made by directly using *hard* model predictions (consider it as a special case of our soft-adversarial strategy with  $\Delta = 0$ ). In addition to controlling the difficulty of the collected examples, it also helps us to reduce the prediction noise from imperfect models and avoid forcing workers to write abnormal events in order to fool the model.

We use two models to provide feedback to the writers, a *future event only* model that focuses primarily on reducing the aforementioned annotation artifacts, and a *premise summary + future event* model that can additionally detect and thus reduce simple negatives that are created as contradictions of the premise. For example, with the premise summary, ‘Howard tells Bernadette that he has a dominant personality’, the negative event ‘Howard will say that he doesn’t have a dominant personality’ is relatively simple as it directly contradicts the premise. Both models are fine-tuned as a sequence classification task from round one’s training data, using a pre-trained RoBERTa-base<sup>4</sup> model. The

objective is to maximize the probability of the positive event being the correct answer. For the future event only model, we only use the future event for classification, ignoring the premise. For the premise summary + future event model, we concatenate the premise summary and the future event text as a single sequence for classification. Note that we use the premise summary as an overall proxy to represent both video and dialogue content to build our adversarial model, considering video and dialogue understanding is still an open research problem in itself.<sup>5</sup> The accuracy of these two models on round one’s Dev. split are 75.34% and 76.68%, respectively. During collection, we randomly pick one model from these two models to provide feedback to users. This is similar to the approach used by Nie et al. (2020) where one model is randomly picked from a set of random seeded models. The difference lies in that we use a set of two models with different inputs (architecture) while Nie et al. (2020) use the same architecture with varying random seeds. This strategy can be seen as constructing a pseudo-ensemble model, which provides diverse adversarial feedback to the crowdworkers and helps avoid annotators exploiting vulnerabilities of a single model (Nie et al., 2020), while reducing server load.<sup>6</sup>

In round two, with our adversarial collection procedure, we collected 7,905 verified examples from

mance but longer response time that affects user experience.

<sup>5</sup>In Appendix Section A.3, we show that an oracle model that uses the premise summary as auxiliary input significantly outperforms our video+dialogue model.

<sup>6</sup>As we only need to run one model instead of multiple models in a standard ensemble approach.

<sup>3</sup>Rewrite for at most twice, in total three trials.

<sup>4</sup>Empirically, RoBERTa-large does not yield better perfor-

4,418 TV show clips and 3,487 YouTube clips. Similar to round one, we split them into 70% training, 15% development, and 15% testing splits.

### 3.4 Adversarial Matching

With adversarial data collection, we are able to collect harder and less-biased examples. However, this approach is not scalable due to its high cost. On average, each verified example in round two costs \$1.70. Inspired by Zellers et al. (2019a), which proposed to use Adversarial Matching to obtain less-biased negatives, we use a similar strategy to create additional examples for our dataset. Given a premise event and its positive event, the goal of adversarial matching is to find a negative from other premise events’ positives, such that the matched negative is very relevant to the premise event (so that they are still hard for machines) and at the same time, not overly similar to the true positive (in case they incidentally become a positive event to the premise). Specifically, we use BERTScore (Zhang et al., 2020) and the recommended RoBERTa-Large model fine-tuned on MNLI (Williams et al., 2018) to calculate similarity score  $S_{sim}(e_i, e_j)$  between two events  $e_i$  and  $e_j$ . For relevance, we use a RoBERTa-base model that takes as input the concatenation of a premise summary  $p_i$  and a future event  $e_j$  and output a relevance score  $S_{rel}(p_i, e_j)$ . This model is trained to distinguish positive events from randomly sampled events. Next, given dataset examples  $\{(p_i, e_i)\}_{i=1}^N$ , we obtain a negative future event for each premise  $p_i$  with maximum-weight bipartite matching (Munkres, 1957; Jonker and Volgenant, 1987) on a weight matrix  $\mathbf{W} \in \mathbb{R}^{N \times N}$ :

$$\begin{aligned} \mathbf{W}_{i,j} &= \lambda(S_{rel}(p_i, e_j) - \alpha S_{sim}(e_i, e_j)), \\ \lambda &= (1 - 0.5 \cdot \mathbb{1}(p_i, e_j)), \end{aligned}$$

where  $\alpha=0.1$  is a hyperparameter that controls the tradeoff between relevance and similarity, the indicator  $\mathbb{1}(p_i, e_j)$  equals 1 if  $p_i$  and  $e_j$  are from different sources (e.g., different TV shows), otherwise 0. Thus,  $\lambda$  serves as a regularization that penalizes  $e_j$  if it is from a different video source than that of  $p_i$  – as  $e_j$  could potentially be an easy negative that can be distinguished from superficial clues such as character names in different shows.

### 3.5 Data Analysis

Table 2 shows the overall statistics of the dataset and data splits details. Each example in our dataset is paired with a premise event clip, with an average

Split	#Videos	Pre. Event	Avg. Sen. Len. (#words)		#Examples
		Avg. Len. (s)	Pre. Sum.	Pos. / Neg.	
Train	7,180	6.1	15.2	11.1 / 11.2	20,142
Dev	1,561	6.2	14.7	11.0 / 11.1	4,392
Test	1,493	6.2	15.4	11.0 / 11.1	4,192
Total	10,234	6.1	15.2	11.1 / 11.2	28,726

Table 2: Statistics by Data Split. *Pre. Event*=Premise Event, a short video with dialogue. *Pre. Sum.*=Premise Summary. *Pos. /Neg.*=Positive/Negative future event.

Domain	Genre	#Shows (#Channels)	#Videos	#Examples
TV show	Sitcom	3	4,117	12,248
	Medical	2	1,558	5,198
	Crime	1	1,072	4,306
YouTube Vlogs	Travel, Food	6	2,406	4,812
	Family, Daily	3	1,081	2,162
Total	-	15	10,234	28,726

Table 3: Data Statistics by Genre.

length of 6.1 seconds. The average length of our positive event (*Pos.*) sentences is very close to that of the negative (*Neg.*) ones (11.1 vs. 11.2), suggesting little bias in sentence length. Our videos are curated from TV shows and YouTube vlogs, across five major categories with diverse topics, i.e., *sitcom*, *medical*, *crime*, *travel-food*, *family-daily*. In Table 3 we show data statistics by genre. Events generally vary by genre. To demonstrate these differences, we show top unique verbs in each genre in Table 4. The top unique verbs in *Crime* genre are usually close to crime and violence, while top unique verbs in *Family, Daily* are usually related to daily activities such as ‘drive’ and ‘wear’. For top unique nouns and additional data analysis (e.g., distribution of examples by reasoning type), please see Appendix Section A.1. For adversarial data collection in round two, the average number of trials is 2.7, i.e., on average the writer has to write their negative event for 2.7 times. For the first trial, 59.21% of the examples are defined as *easy* by our system, i.e., the positive event has a much larger probability of happening than the negative event. With rewriting, only 31.22% of the examples remain *easy*. Moreover, in Table 7 row 1, when trained on our final dataset, we show that our future event only baseline gets much lower performance on the round two subset than that of round one (59.62% vs. 74.20%), showing round two examples are less-biased.

Genre	Top Unique Verbs
Sitcom	change, offer, hear, should, accept, yell, hang, join, apologize, shut, shout, realize
Medical	die, treat, cry, yell, smile, proceed, examine, approach, argue, save, admit, rush
Crime	kill, shoot, point, question, toss, hang, remove, catch, lie, deny, investigate,
Travel, Food	taste, add, pour, dip, cook, describe, cut, order, serve, stir, prepare, enjoy, buy
Family, Daily	drive, jump, wear, point, smile, touch, climb, dress, set, swim, hide, lay, blow

Table 4: Top unique verbs in each genre.

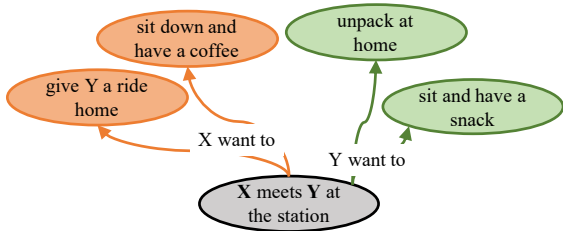


Figure 3: An ATOMIC (Sap et al., 2019) example.

## 4 Method

Given a video with dialogue text, and two future event candidates  $\{e_i\}, i \in \{1, 2\}$ , our goal is to predict which future event is more likely to happen. In the following, we introduce our baseline approach (see model overview in Figure 4) for this new task.

**Video Encoding.** We encode each video using appearance and motion features at 1 FPS. For appearance, we extract 2048D feature vectors from the ImageNet (Deng et al., 2009) pre-trained ResNet-152 (He et al., 2016). For motion, we extract 2048D feature vectors from the Kinetics (Carreira and Zisserman, 2017) pre-trained ResNeXt-101 (Hara et al., 2018). These features have shown to perform well in several video and language tasks (Miech et al., 2019). We perform L2-normalization and concatenate the features as the video representation. We project these representations into a lower dimension space and add a trainable positional encoding (Devlin et al., 2019) to them. We then use a *transformer encoder* (Vaswani et al., 2017) to further encode the resulting representation, denoted as  $E^v \in \mathbb{R}^{T \times d}$ .

**Text Encoding.** For text, we use the contextualized text features from the RoBERTa-base (Liu et al., 2019). We first fine-tune the pre-trained RoBERTa with commonsense knowledge extracted from the ATOMIC dataset (Sap et al., 2019) (see

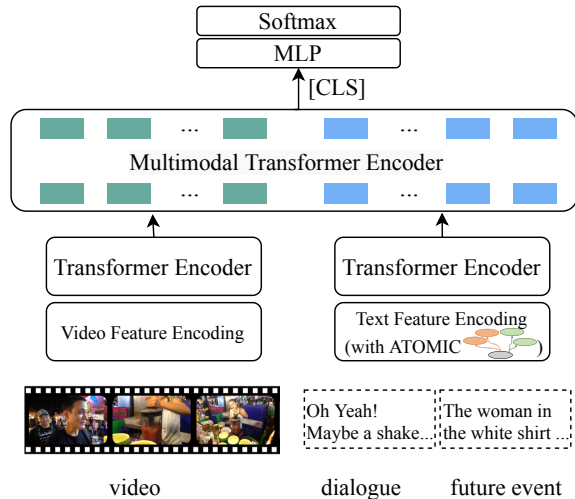


Figure 4: Model overview. We first separately encode video and text, and then use a multimodal transformer encoder to encode information from both modalities. Please see text for details.

details in the paragraph below) and then use the resulting model for feature encoding. Note that this model is end-to-end trainable during training. We concatenate dialogue and future event candidate as input to the transformer layers, special tokens such as [CLS] (Devlin et al., 2019) are also added in this process. We use the extracted token embeddings from the last layer, denoted as  $E_i^t \in \mathbb{R}^{L_i \times d}, i \in \{1, 2\}$ , where  $L_i$  is the sequence length (#tokens, including added special tokens). Similar to how we encode video, the resulting text representation is further encoded using a transformer encoder. Without ambiguity, we use the same notation to denote the outputs as  $E_i^t \in \mathbb{R}^{L_i \times d}, i \in \{1, 2\}$ .

**Commonsense-based Text Representations.** Addressing our challenging future event prediction task requires general world knowledge that is beyond basic visual and language semantic understanding. Thus, we propose to inject the commonsense from the ATOMIC dataset (Sap et al., 2019) into our framework in a simple way. ATOMIC contains events with if-then inferences, e.g., *if X meets Y at the station, then X want to give Y a ride home* (see example in Figure 3). We extract 406K event inferences from the dataset, and replace the person tokens X and Y with the names from our dataset (Mitra et al., 2019). We then use the extracted event inference sentences to finetune the pre-trained RoBERTa-base model. The fine-tuned model is then used to encode our text inputs.

Model	Accuracy (%)
chance	50.00
future only	58.09
video + future	59.03
dialogue + future	66.63
video + dialogue + future	67.46
human (dialogue + future)	76.25
human (video + dialogue + future)	90.50

Table 5: Results on VLEP Test split.

## Multimodal Encoding and Event Classification.

To obtain the joint multimodal representation, we concatenate encoded video  $E^v$  and text  $E^t$  and use a transformer encoder to encode the concatenated representations. This encoder allows information exchange between the two modalities. We use the representation from the [CLS] token as the joint representation of video, dialogue and future event  $e_i$ , denoted as  $g_i \in \mathbb{R}^d, i \in \{1, 2\}$ . We gather the joint representation vectors for all future event candidates and pass them through a two-layer MLP with a softmax layer for classification. We train the model using cross-entropy loss that maximizes the scores for the more-likely future events.

## 5 Experiments

### 5.1 Implementation Details

Our models are implemented in PyTorch (Paszke et al., 2017). To speed up training, we use NVIDIA Apex for mixed precision training. We set the hidden size  $d$  to be 768 and use a single transformer layer for all our transformer encoders. We use Adam (Kingma and Ba, 2015) optimizer with  $\beta_1=0.9, \beta_2=0.999$ . Since our model has a pre-trained component (RoBERTa), we use a two-phase training strategy. Specifically, we first freeze RoBERTa’s weights up to the second last layer and then pre-train the rest of model for 3 epochs with initial learning rate of  $1e-4$ , learning rate warmup over the first 10% of the steps and linear decay the learning rate to 0. We then unfreeze all the weights and finetune the whole model for 3 epochs with learning rate  $5e-5$  and linearly decay the learning rate to 0. We train the model on a single RTX 2080Ti GPU with batch size 16. We report multiple-choice question answering accuracy.

### 5.2 Results

**Are video and dialogue modalities useful?** Table 5 shows the results with different input context. The model using future event text only as the input

Model	Accuracy (%)
video + dialogue + future	67.46
- ATOMIC fine-tuning	66.96

Table 6: Effect of ATOMIC fine-tuning.

Model	Adv. Matching	Human-Annotated		Overall
	(50%)	R1 (22%)	R2 (28%)	
future only	50.00	74.20	59.62	58.09
video + future	54.34	69.21	59.19	59.03
dialogue + future	67.60	70.70	61.53	66.63
video + dialogue + future	68.37	70.59	63.26	67.46

Table 7: Performance breakdown by data collection method.

achieves 58.09% accuracy, which is higher than random chance (50%), suggesting there exists slight bias even with our deliberate adversarial collection and matching but is tolerable. Adding video or dialogue as additional input improves the accuracy to 59.03% and 66.63%, respectively. The best performance is achieved when using both video and dialogue, with an accuracy of 67.46%. In Appendix Section A.3, we also present an oracle model with premise summary as auxiliary input.

**Human Performance.** To obtain human performance, we randomly sampled 400 examples from our test set. We present a premise event (a video with dialogue subtitles or dialogue subtitles only) and its two corresponding future events to a new set of AMT workers and ask them to select which one is more likely to happen after the premise. Each example is answered by 10 different workers to reduce crowdworker variance (Rajpurkar et al., 2018). The final answer is selected by majority vote. Table 5 shows the results. We observe that human performance without video (i.e., only dialogue+future) is 76.25%, while showing the video improves the performance to 90.5%. which indicates video information is important for getting the correct answer. Compared with the best model result (67.46%), there is still a large useful gap (23%) for future community work on our challenging task of multimodal event prediction.

**Does commonsense knowledge help?** In Table 6, we show a model variant that uses text features without ATOMIC sentences for fine-tuning. We see this variant has a lower accuracy (66.96%) compared with the fine-tuned accuracy (67.46%).

**Impact of Data Collection Method.** Table 7 shows the model performance breakdown by differ-



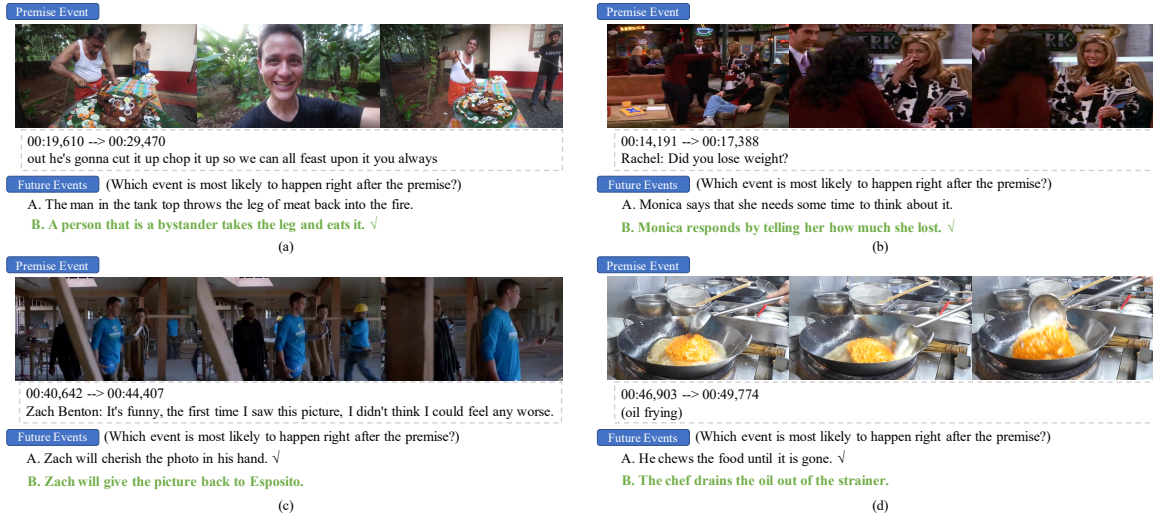


Figure 5: Prediction examples from our best model. Top row shows correct predictions, bottom row shows failure cases. Left column shows human annotated examples, right column shows adversarial matching examples. Ground truth answers are in bold and green, model predictions are indicated by ✓.

word	throws	face	leave	without
PMI (R1)	1.38	1.28	1.25	1.23
PMI (R2)	0.83	0.67	0.66	0.81

Table 8: Top words by PMI in standard collection (R1) and their values in adversarial collection (R2). The values are calculated using PMI(word, less-likely).

ent collection methods. For human-annotated data, we show performance on round one (*R1*, standard data collection) and round two (*R2*, adversarial data collection). First, we observe that the accuracy of the *future only* model matches chance on adversarial matching data while being higher on human-annotated data. The main reason is the matched data has less artificial biases than human-annotated ones. Second, for human-annotated data, across all models, we see the performance on round two subset is significantly lower than that of round one, which demonstrates the effectiveness of using our adversarial collection procedure.

Gururangan et al. (2018) shows lexical choice is a strong indicator of the inference class in NLI. To check how our adversarial collection affects the use of words, we use pointwise mutual information (PMI) as in Gururangan et al. (2018). In Table 8 we show top words that are associated with negative class (less-likely event) in standard collection versus their values in our adversarial collection process. We find that the PMI values of these top negative words (e.g., ‘throws’, ‘without’, that frequently occur in negative less-likely events) in standard collection clearly drop in adversarial collection, e.g.,

‘throws’ drops from 1.38 to 0.83, making it less indicative of the negative.

**Qualitative Examples.** We show 4 prediction examples using our best model (video + dialogue + future) in Figure 5. Top row shows two correct prediction examples, where our model is able to predict basic human intention and reaction. Bottom row shows two incorrect predictions, where wrong predictions are mainly caused by the lack of commonsense. For example, to correctly pick the more likely event in Figure 5(c), the model needs to understand that the ‘photo’ is an evidence of a police investigation. Figure 5(d) shows an example that requires the model to infer that the food is not ready yet. More examples are presented in Appendix Section A.4.

## 6 Conclusion

We introduce a new task, Video-and-Language Event Prediction (VLEP) - given a video with aligned dialogue, and two future events, an AI system is required to predict which event is more likely to happen. To support this task, VLEP dataset is collected. We present a strong transformer-based baseline that incorporates information from video, dialogue, and commonsense knowledge, each of which is necessary for this challenging task.

## Acknowledgements

We thank the reviewers for their helpful feedback. This research is supported by NSF Award

#1562098, DARPA MCS Grant #N66001-19-2-4031, DARPA KAIROS Grant #FA8750-19-2-1004, ARO-YIP Award #W911NF-18-1-0336, and Google Focused Research Award. The views contained in this article are those of the authors and not of the funding agency.

## References

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *ICCV*.
- Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*.
- Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. 2015. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*.
- Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*.
- Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the ninth workshop on statistical machine translation*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.
- Dave Epstein, Boyuan Chen, and Carl Vondrick. 2020. Oops! predicting unintentional action in video. In *CVPR*.
- Zhiyuan Fang, Tejas Gokhale, Pratyay Banerjee, Chitta Baral, and Yezhou Yang. 2020. Video2commonsense: Generating commonsense descriptions to enrich video captioning. In *EMNLP*.
- David F Fouhey, Wei-cheng Kuo, Alexei A Efros, and Jitendra Malik. 2018. From lifestyle vlogs to everyday interactions. In *CVPR*.
- Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. 2015. Recurrent network models for human dynamics. In *ICCV*.
- Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. 2017a. Tall: Temporal activity localization via language query. In *ICCV*.
- Jiyang Gao, Zhenheng Yang, and Ram Nevatia. 2017b. Red: Reinforced encoder-decoder networks for action anticipation. In *BMVC*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R Bowman, and Noah A Smith. 2018. Annotation artifacts in natural language inference data. In *NAACL*.
- Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. 2018. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *CVPR*.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*.
- Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. 2017. Localizing moments in video with natural language. In *ICCV*.
- Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. 2018. Multimodal explanations: Justifying decisions and pointing to the evidence. In *CVPR*.
- Oana Ignat, Laura Burdick, Jia Deng, and Rada Mihalcea. 2019. Identifying visible actions in lifestyle vlogs. *ACL*.
- Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. 2017. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *CVPR*.
- Roy Jonker and Anton Volgenant. 1987. A shortest augmenting path algorithm for dense and sparse linear assignment problems. *Computing*.
- Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. 2014. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*.
- Kyung-Min Kim, Min-Oh Heo, Seong-Ho Choi, and Byoung-Tak Zhang. 2017. Deepstory: Video story qa by deep embedded memory networks. In *IJCAI*.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.
- Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara L Berg, and Mohit Bansal. 2020a. Mart: Memory-augmented recurrent transformer for coherent video paragraph captioning. In *ACL*.
- Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L. Berg. 2018. Tvqa: Localized, compositional video question answering. In *EMNLP*.

- Jie Lei, Licheng Yu, Tamara L Berg, and Mohit Bansal. 2020b. Tvqa+: Spatio-temporal grounding for video question answering. In *ACL*.
- Jie Lei, Licheng Yu, Tamara L Berg, and Mohit Bansal. 2020c. Tvr: A large-scale dataset for video-subtitle moment retrieval. In *ECCV*.
- Xiaodan Liang, Lisa Lee, Wei Dai, and Eric P Xing. 2017. Dual motion gan for future-flow embedded video prediction. In *ICCV*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *ACL*.
- Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*.
- Jingzhou Liu, Wenhui Chen, Yu Cheng, Zhe Gan, Licheng Yu, Yiming Yang, and Jingjing Liu. 2020. Violin: A large-scale dataset for video-and-language inference. In *CVPR*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Wei Mao, Miaomiao Liu, Mathieu Salzmann, and Hongdong Li. 2019. Learning trajectory dependencies for human motion prediction. In *ICCV*.
- R Thomas McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *ACL*.
- Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. 2019. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *ICCV*.
- Arindam Mitra, Pratyay Banerjee, Kuntal Kumar Pal, Swaroop Mishra, and Chitta Baral. 2019. Exploring ways to incorporate additional knowledge to improve natural language commonsense question answering. *ArXiv*, abs/1909.08855.
- James Munkres. 1957. Algorithms for the assignment and transportation problems. *Journal of the society for industrial and applied mathematics*.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial nli: A new benchmark for natural language understanding. In *ACL*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *ACL*.
- Jae Sung Park, Chandra Bhagavatula, Roozbeh Motlaghi, Ali Farhadi, and Yejin Choi. 2020. Visual commonsense graphs: Reasoning about the dynamic context of a still image. In *ECCV*.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*.
- Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for squad. In *ACL*.
- Anna Rohrbach, Atousa Torabi, Marcus Rohrbach, Niket Tandon, Christopher Pal, Hugo Larochelle, Aaron Courville, and Bernt Schiele. 2017. Movie description. *IJCV*.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for if-then reasoning. In *AAAI*.
- Yuge Shi, Basura Fernando, and Richard Hartley. 2018. Action anticipation with rbf kernelized feature mapping rnn. In *ECCV*.
- Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2016. Movieqa: Understanding stories in movies through question-answering. In *CVPR*.
- Masatoshi Tsuchiya. 2018. Performance impact caused by hidden bias of training data for recognizing textual entailment. In *LREC*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015a. Cider: Consensus-based image description evaluation. In *CVPR*.
- Ramakrishna Vedantam, Xiao Lin, Tanmay Batra, C Lawrence Zitnick, and Devi Parikh. 2015b. Learning common sense through visual abstraction. In *ICCV*.
- Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. 2016a. Anticipating visual representations from unlabeled video. In *CVPR*.
- Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. 2016b. Generating videos with scene dynamics. In *NeurIPS*.

Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. 2019. Vatex: A large-scale, high-quality multilingual dataset for video-and-language research. In *ICCV*.

Adina Williams, Nikita Nangia, and Samuel R Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *NAACL*.

Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*.

Amir Zadeh, Michael Chan, Paul Pu Liang, Edmund Tong, and Louis-Philippe Morency. 2019. Social-iq: A question answering benchmark for artificial social intelligence. In *CVPR*.

Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019a. From recognition to cognition: Visual commonsense reasoning. In *CVPR*.

Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. Swag: A large-scale adversarial dataset for grounded commonsense inference. In *EMNLP*.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019b. Hellaswag: Can a machine really finish your sentence? In *ACL*.

Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2016. Yin and yang: Balancing and answering binary visual questions. In *CVPR*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *ICLR*.

## A Appendices

### A.1 Additional Data Analysis

Our videos are curated from two sources, TV shows and YouTube lifestyle vlogs, across five major categories, i.e., *sitcom*, *medical*, *crime*, *travel-food*, *family-daily*. Events generally vary by genre. One way to show the difference is by checking the top unique nouns in each genre. To obtain the top unique nouns, we first tokenize and lemmatize the future event sentences. Each resulting token is also tagged with a part-of-speech tag. Next, for each genre, we take the top unique nouns as the ones among the most frequent 100 nouns from one genre but do not appear in those from the other genres combined. We show the top unique nouns in each genre in Table 9. Interestingly, the top unique nouns in *crime* genre are closer to crime and violence, while in *family-daily*, the top unique nouns are relatively more family relevant.

Genre	Top Unique Nouns
Sitcom	apartment, group, couch, bottle, game, date, joke, kitchen, story, wine, seat, hug
Medical	patient, doctor, office, surgery, parent, elevator, hospital, nurse, team, case, cane
Crime	gun, picture, photo, paper, information, evidence, police, case, suspect, ground
Travel, Food	host, meat, bite, plate, bowl, chef, piece, sauce, fish, dish, soup, noodle, spoon
Family, Daily	kid, dad, child, dog, son, toy, father, daughter, family, wife, video, candy, hair

Table 9: Top unique nouns in each genre.

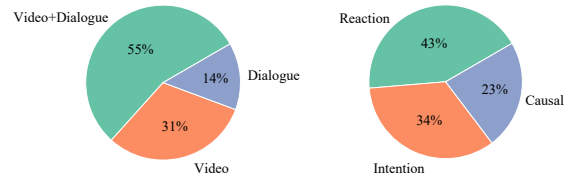


Figure 6: Distribution of examples by premise understanding type (left) and by reasoning type (right).

Figure 6 (left) shows the distribution of examples by premise understanding type, i.e., what modalities are needed to understand the premise event. Most of the premise events require both video and dialogue understanding. Figure 6 (right) shows the distribution of examples by commonsense reasoning type. We categorize commonsense reasoning into three types by examining the relation between the premise event and the positive future event: (1) intention, e.g., if *X brings two cups of coffee*, then *X (intends to) give Y a cup of coffee*. (2) reaction, e.g., if *X hands Y a form and describes a procedure*, then *Y signs the form and hands it back*. (3) causal, e.g., if *X says they hit a bump*, then *X gets unbalanced and falls off the boat*. The two distributions are obtained by manually annotating 100 randomly sampled examples from VLEP Dev. split.

Next, we show the distribution of premise event length and premise summary length in Figure 7 and Figure 8, respectively. In addition, we also show the distribution of positive future event length and negative event length in Figure 9 and Figure 10.

### A.2 Additional Data Collection Details

We hire workers from Amazon Mechanical Turk (AMT) to annotate our data. To ensure our data quality, we only allow workers from English-speaking countries to participate in our task. We require workers to have at least 500 HITs approved

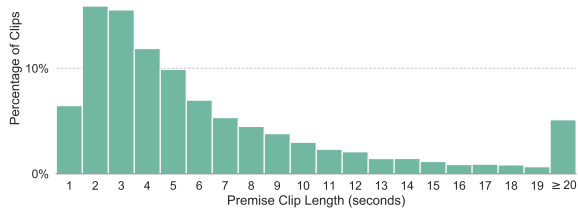


Figure 7: Distribution of premise event length.

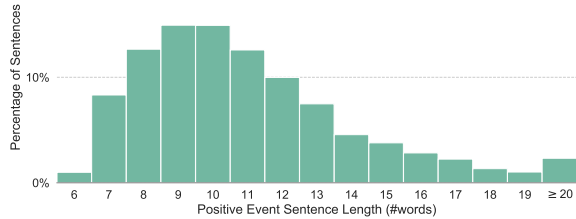


Figure 9: Distribution of positive future event length.

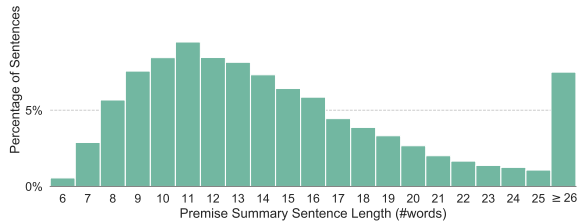


Figure 8: Distribution of premise summary length.

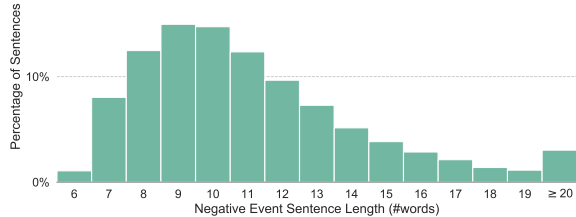


Figure 10: Distribution of negative future event length.

with an approval rate of 95%. Furthermore, we design a qualification test with 10 multiple-choice questions to ensure that workers well understand our annotation requirement. We show an example question from our qualification test in Figure 11. The workers have to correctly answer at least 7 questions to pass the test. In total, 518 workers participated the test, with a pass rate of 56%. During data collection, we set up an automatic tool to check if all required annotations have been performed. We also manually review the submitted results and provide prompt feedback to them, encouraging better annotation.

Our data collection instructions and interface for round two (adversarial data collection) are shown in Figure 12 and Figure 13, respectively. Round one collection details are similar to that of the round two, except that we do not require workers to fool our basic models (*robot*). In our annotation process, the actual future events in the videos are not hidden from the workers to ease the collection. The workers can either write the actual future event as the more likely event, or they can hypothesize one when the actual future event in the given video is surprising/rare (such as some events in sitcoms). To ensure the quality of the examples, we conduct a strict filtering step in which each example is verified by three extra workers (verifiers) and we only accept examples where at least three out of four (one writer + three verifiers) reach an agreement, as Hendricks et al. (2017); Nie et al. (2020).

Model	Accuracy (%)
video + dialogue + future	67.46
+ premise summary (oracle)	75.64

Table 10: Oracle performance with premise summary.

### A.3 More Results

**Oracle Premise Results.** As an oracle test, we apply the collected premise summary as an auxiliary input to the model, removing certain obstacles of video-dialogue understanding in our baseline model. We show this oracle model performance in Table 10. Our model with premise summary (oracle) achieves 75.64%, which is significantly higher than the one without it (67.46%), indicating the desire for better video-dialogue understanding.

**Future Event Generation Results.** Given the videos, we can also set up an alternative task of using a captioning-style model to generate future event descriptions. We use the MultiModal Transformer from Lei et al. (2020c) as our baseline for this task. This model uses a standard transformer encoder-decoder architecture for caption generation. The video embeddings and dialogue embeddings are concatenated as inputs (Lei et al., 2020a) to the transformer encoder. We use the default model and training configurations from Lei et al. (2020c). With this system, we evaluate generation performance with video and dialogue as inputs. Our video+dialogue model has CIDEr-D (Vedan-

tam et al., 2015a): 19.57, BLEU@4 (Papineni et al., 2002): 1.80, Rouge-L (Lin, 2004): 16.42, and METEOR (Denkowski and Lavie, 2014): 7.58. Note that we only use this generation task to demonstrate that it is possible to generate future event sentences from videos. This may not be as suitable as our default multiple choice setup to serve as a benchmark, since generation is known to be more difficult to evaluate (Liu et al., 2016). Besides, it also requires multiple references (Vedantam et al., 2015a) to be more accurate. Therefore, we recommend future work to use human evaluation if you pursue a generation-based setup on our dataset.

#### **A.4 More Qualitative Examples**

We show more correct and incorrect predictions from our best model (video + dialogue + future) in Figure 14 and Figure 15, respectively.



Premise Summary A panini, billowing smoke, descends in front of Raj, Howard and Leonard.



check premise clip

**Future event:** (what might happen right after the premise clip). **Rationale:** (why the event is more or less likely, be concise.)

more-likely	Someone picks up the panini and takes it off the surface.	rationale	After cooking, people will take the food out.
less-likely	The panini sets fire to the surface and the guys panic.	rationale	There is no visible fire and the surface doesn't seem to be combustible.

Question: Given only the premise event/clip, is this example correct? Why?

- A. Yes!
- B. No! Timestamps are not correct, they cover irrelevant clips.
- C. No! The more-likely event is less likely to happen than the less-likely event.
- D. No! The future events are not related to the premise event (clip).

Figure 11: Example question from our qualification test. Workers have to correctly answer 7 out of 10 questions in the test to participate in our annotation task.

**Instructions**

Find out premise events in a video that may lead to some possible future events to happen. Write down both the premise and the future events.

**Steps:**


1. Watch the video by clicking the video player, find out a **premise event**, write a summary for it.
2. Drag slider handles to locate the clip that **precisely shows only the premise event**.
3. Write two **future events**:
  - **more-likely event (>50% chance)**: an event that is very likely to happen after the premise.
  - **less-likely event (<50% chance)**: a related event that is less likely to happen. It cannot be a random or impossible event.
4. For each future event, write a **rationale**, indicating why the event is likely or less-likely.

**FAQ (MUST READ):**

- **What is an event?**  
Topics/things discussed in dialogues (subtitles) and activities/actions of people/animals etc., in the videos. A single event does not necessarily mean a single action, it could also be a combination of multiple concurring actions.
- **What are the general requirements for an event description?**  
(1) It should be written in standard English and contains at least 6 words. (2) When writing dialogue related events, do not copy the text in the dialogue word-by-word, please paraphrase the original text with your own words.
- **Special requirements for less-likely events.**  
**Goal: Write a tricky less-likely event in a way such that humans can distinguish it from the more-likely event, but the robot will get fooled.** We have a **smart robot** that looks at the videos and **learns from your previous events**. It knows that one event is less likely to happen if it sees the events that are similar to this event are all less-likely events. Thus please don't write obvious repeated events.
- **How do I mention people in the descriptions?**  
For videos with names in the subtitles, use the **names** of the characters you are referring to in your description. When names are not present, please provide a unique description of the person, e.g., 'the man in red hat'.
- **What are the requirements for an event rationale?**  
(1) It should be written in standard English and contains at least 6 words. (2) **Be concise**, only write the direct reason that explains why the corresponding event is more or less likely. For example, 'House starts to do jumping jacks is less likely because House has a disability with his leg.' can be simplified as 'House has a disability with his leg.'
- **What are the requirements for locating a premise event clip?**  
Make the START/END timestamps of the clip tight, i.e., **cover the premise event completely and do not cover irrelevant parts**.
- **Which events do you prefer?**  
There is no specific event we would like you to write. **But please try your to write different events for different videos, i.e., do **\*\*not\*\*** repeat yourself or use similar patterns to write the events.**

Figure 12: Annotation instructions for round two (adversarial data collection).

**HIT starts here**



Premise Summary

00:00  02:00 check premise clip

**Future events:** (what might happen right after the premise clip). **Rationale:** (why the event is more or less likely, be concise.).

<b>more likely</b>	What is most likely to happen right after the premise?	<b>rationale</b>	Why is it likely? Please be concise.
<b>less likely</b>	What is less likely to happen than the likely one above? It cannot be a random or impossible event.	<b>rationale</b>	Why is it less likely? Please be concise.

Write the less-likely event in the yellow box above, this box will show you robot answers.

**Goal: Write a tricky less-likely event such that humans can distinguish that it is less likely to happen than the more-likely event, but the robot will get fooled.** This smart robot looks at the videos and **learns from your previous events**. It gets to know that one event is less likely to happen if it sees the events that are similar to this event are mostly less-likely events. Thus please write different events, avoid writing obvious repeated events in your different submissions. Click the following button before submitting to see further instructions. Get Robot Answer

Figure 13: Annotation interface for round two (adversarial data collection).



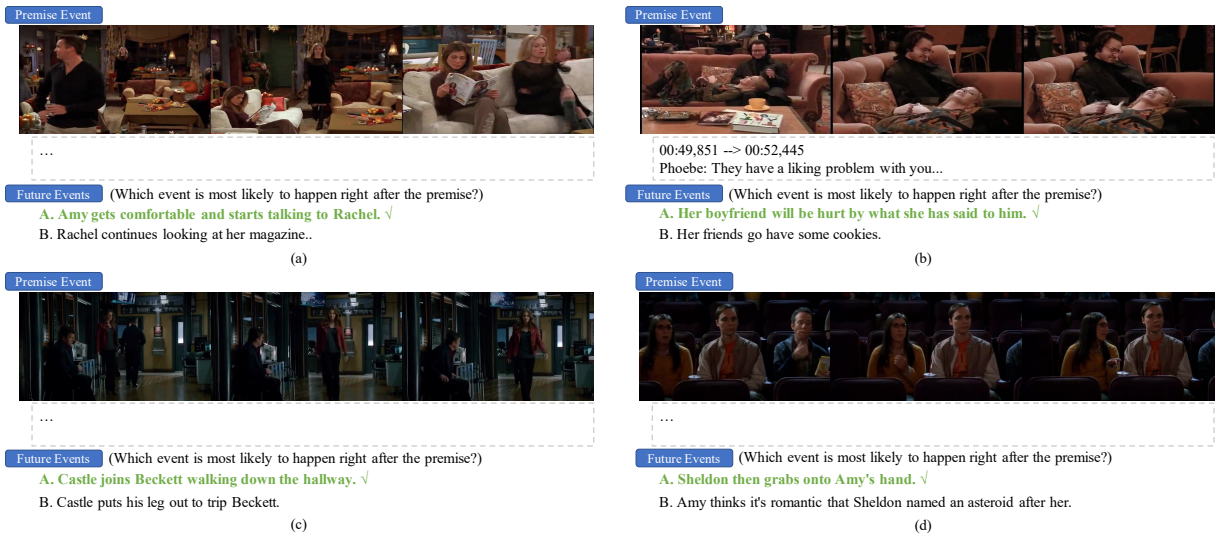


Figure 14: Correct prediction examples from our best model. Left column shows human annotated examples, right column shows adversarial matching examples. Ground truth answers are in bold and green, model predictions are indicated by ✓.

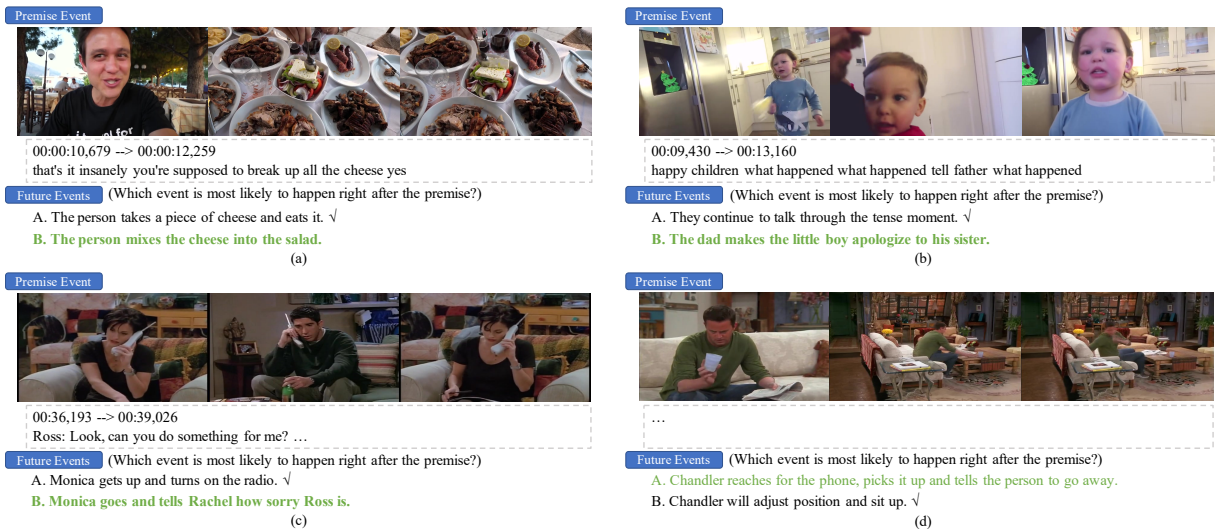


Figure 15: Failure examples from our best model. Left column shows human annotated examples, right column shows adversarial matching examples. Ground truth answers are in bold and green, model predictions are indicated by ✓.