

Synopses of Movie Narratives: a Video-Language Dataset for Story Understanding

Anonymous ACL submission

Abstract

Despite recent advances of AI, story understanding remains an open and under-investigated problem. We collect, preprocess, and publicly release a video-language story dataset, Synopses of Movie Narratives (SYMON), containing 5,193 video summaries of popular movies and TV series. SYMON captures naturalistic storytelling videos for human audience made by human creators, and has higher story coverage and more frequent mental-state references than similar video-language story datasets. Differing from most existing video-text datasets, SYMON features large semantic gaps between the visual and the textual modalities due to the prevalence of reporting bias and mental state descriptions. We establish benchmarks on video-text retrieval and zero-shot alignment on movie summary videos. With SYMON, we hope to lay the groundwork for progress in multimodal story understanding.

1 Introduction

Stories are complex artifacts that succinctly encode the human experience. The understanding of story content involves high-level semantic concepts such as character motivations and intentions (Emelin et al., 2020; Rashkin et al., 2018), events structures (Chambers and Jurafsky, 2008; Li et al., 2013; Pichotta and Mooney, 2016; Ferraro and Van Durme, 2016; Martin et al., 2018; Wang et al., 2021; Caselli et al., 2021), as well as social relationships among story characters (Elson et al., 2010; Chaturvedi et al., 2016; Kim and Klinger, 2019). To this day, understanding of story semantics remains an open and under-investigated problem.

The recent emergence of user-generated, “a movie in X minutes” videos offers a rich source of naturalistic storytelling videos. These videos usually select clips that depict key story events from a movie or a TV series. The narrator recounts the story alongside the video. These videos provide



Figure 1: An example video with narration text from SYMON. The video has been automatically segmented into three scenes. We show the boundary timestamps.

condensed yet complete storylines that are carefully assembled for human viewers by human creators.

We identify, collect, preprocess, and publicly release a video-language story dataset, named Synopses of Movie Narratives (SYMON). The dataset includes 5,193 user-generated video summaries of popular movies and TV series for a total length of 869 hours. For 857 movies, multiple summary videos are available, which may be used as references for generation or summarization. In Figure 1, we show an example video and text description from SYMON. We empirically verify SYMON as the prototypical story dataset, as it has higher coverage of plotlines and more frequent mental-state references than several similar video-language story datasets.

However, the nature of storytelling poses unique obstacles for computational understanding due to the semantic divergence between the video and text. First, in the phenomenon known as reporting bias (Gordon and Van Durme, 2013), human narrators tend to avoid stating the obvious. For example, in Figure 1, the video shows Harry Pot-

065 ter lying on the floor, while the narrator states "...
066 knocking him unconscious". To recognize that lying
067 on the floor is a consequence of being knocked
068 unconscious requires event-level cause-and-effect
069 reasoning, which may prove difficult for today's
070 AI (Sap et al., 2019). Second, the story texts con-
071 tain frequent mentions of story characters' mental
072 states (§5.2), which may not be easily recogniz-
073 able from video. This contrasts with crowdsourced
074 datasets like Charades (Sigurdsson et al., 2016)
075 where humans are asked to follow textual instruc-
076 tions, or LSMDC (Rohrbach et al., 2017) where
077 the narration meticulously describes the imagery
078 for audience with visual impairment.

079 To examine the cross-modality semantic gap, we
080 design a simple task that temporally orders two
081 video segments. A large pretrained UniVL model
082 (Luo et al., 2020) demonstrates mediocre perfor-
083 mance and limited utilization of textual informa-
084 tion, highlighting the challenge posed by SYMON.

085 As benchmarks for future research, we estab-
086 lish baselines for text-to-video and video-to-text
087 retrieval on SYMON and a zero-shot video-text
088 alignment baseline using the YMS dataset as test.
089 Together, the weakly supervised SYMON and the
090 fully annotated YMS form a complete benchmark,
091 serving as a new challenge for the multimodal re-
092 search community.

093 Our contributions are three-fold:

- 094 • We collect, preprocess, and publish a large-
095 scale movie summary dataset, which can sup-
096 port various multimodal tasks such as re-
097 trieval, captioning, and summarization.
- 098 • We perform extensive experiments to quantify
099 the characteristics of SYMON, including its
100 coverage of major plotlines, the amount of
101 mental-state descriptions, and the semantic
102 divergence between text and video.
- 103 • To facilitate future research, we establish base-
104 lines for text-video retrieval on SYMON and
105 zero-shot transfer to the YouTube Movie Sum-
106 mary dataset (YMS) (Dogan et al., 2018).

107 2 Related Work

108 **Datasets for Event and Story Understanding.**
109 Events and story structures are closely related
110 (Caselli et al., 2021). Existing datasets provided
111 annotations for the temporal aspects, such as
112 temporal precedence and duration (UzZaman et al.,

113 2013; Chambers et al., 2014; Ning et al., 2020;
114 Zhou et al., 2021; Vashishtha et al., 2019, 2020),
115 and causal relations between events (O’Gorman
116 et al., 2016; Roemmele et al., 2011).

117 Several datasets explore individual components
118 of stories, including sentence ordering (Gangal
119 et al., 2021), social norms and moral consequences
120 (Emelin et al., 2020), plausible antecedent (Bha-
121 gavatula et al., 2020), intentions and effects on
122 mental states (Rashkin et al., 2018), high-level
123 story structures (Ouyang and McKeown, 2015; Li
124 et al., 2018), and story character descriptions (Brah-
125 man et al., 2021). Sap et al. (2019) consider rela-
126 tions between events, persona, and mental states.
127 Some datasets aim at summarization for screen-
128 plays or conversation transcripts (Gorinski and La-
129 pata, 2015; Papalampidi et al., 2020; Chen et al.,
130 2021). Notably, Sadhu et al. (2021) annotate event
131 relations from video.

132 Researchers also develop general-purpose QA
133 datasets conditioned on comprehension of story
134 texts, such as MCTest (Richardson et al., 2013),
135 NarrativeQA (Kočíský et al., 2018), and FriendsQA
136 (Yang and Choi, 2019). Multimodal counterparts
137 like MovieQA (Tapaswi et al., 2016), TVQA (Lei
138 et al., 2018), and Pororo (Kim et al., 2017) are
139 available. However, not every question in the QA
140 datasets requires in-depth narrative understanding.

141 **Video-Text Movie Story Datasets.** A number
142 of datasets supply story content extracted from
143 movies. The Large-Scale Movie Description Chal-
144 lenge (LSMDC) (Rohrbach et al., 2017) combined
145 the efforts of MPII-MD (Rohrbach et al., 2015)
146 and M-VAD (Torabi et al., 2015) and provide de-
147 tailed language descriptions initially intended for
148 the visual impaired. Although these descriptions
149 are highly accurate, they may not be representative
150 of real-world storytelling.

151 YouTube Movie Summary (YMS) (Dogan et al.,
152 2018) contains 94 YouTube movie summary videos
153 with human-narrated storylines. The Condensed
154 Movies Dataset (CMD) (Bain et al., 2020) gath-
155 ers 7 to 11 key clips from each movie with one-
156 sentence descriptions for each clip. Pororo (Kim
157 et al., 2017) captures 20-minute cartoon episodes,
158 in-show conversations, and human-written descrip-
159 tions. MovieNet (Huang et al., 2020) annotate 2000
160 hours of movies with extensive annotations and
161 aligned movies scripts. However, due to copyright,
162 legal clearance for the video release is still pending
163 at the time of writing.

	Video hours	#Videos (#Clips)	#Sent	Vocab.
CMD	1,270	3,605 (33,976)	35,681	15,272
MovieNet (video release pending)	2,000	1,100		
LSMDC	147	200 (128,085)	128,118	22,500
Pororo	20.5	171 (16,066)	43,394	
MovieGraph	94.0	51 (7,637)	20,849	
SYMÓN (Ours)	869	5,193	683,611	40,116

Table 1: Comparison of video description datasets with story content.

Other types of video annotations have been explored, including semantic roles and event relations (Sadhu et al., 2021), character relationships and types of speech (Wu and Krahenbuhl, 2021), and movie graphs (Vicol et al., 2018).

In this work, we collect a large-scale, readily available, multi-reference dataset of human-curated movie summaries, named SYMÓN. The dataset can be leveraged for various story understanding and generation tasks such as sequential text localization, story generation from video, and movie summarization. To our knowledge, SYMÓN is the largest dataset for short naturalistic storytelling videos.

3 Dataset Collection and Statistics

We apply the following procedure for data collection. First, we manually identify relevant YouTube channels by searching with keywords such as “movie summary”, “movie recap”, and “movie shortened”. We download all videos from the identified channels and accompanying subtitles, which may be written by humans or automatically generated by YouTube. Videos without subtitles are excluded. Finally, we perform rule-based extraction of movie names from metadata and subtitles and discard videos that are not movie summaries.

This yields a total of 5,193 videos with an average length of 9.5 minutes and a total length of 869 hours. On average, the narration in one video contains 1,717 words or 131 sentences. The overall vocabulary size is 40,116. SYMÓN contains summaries for 2,440 movies and TV series, of which 857 have more than 1 summary. The most popular TV series, *The Walking Dead*, has 84 summaries. On average, one movie or TV series in the 857 has 4.21 summaries. Compared to existing datasets (see Table 1) SYMÓN is one of the largest movie

narrative datasets with most diverse vocabulary. In addition, SYMÓN has more complete coverage of story events than LSMDC and CMD (§5.1).

4 Preprocessing

Subtitle Masking. Some videos have subtitles embedded in the video. In tasks like text-to-video retrieval, the embedded subtitles may become a shortcut feature, causing networks to learn only optical character recognition.

To eliminate shortcuts, we locate embedded subtitles and mask them out. For efficiency, we randomly sample 100 frames from each video and apply an accurate text detection technique (Baek et al., 2019). Observing that the subtitles are almost always at the same location in a single video, we take the minimum bounding box that can cover all embedded subtitles in all 100 frames as the masked region; we set all pixels in the region to black.

Punctuation Restoration. We acquire subtitle texts from YouTube directly. Sometimes the texts are the result of automatic speech recognition, which cannot recognize punctuation. To fix this, for every unpunctuated narration text, we generate punctuation with (Alam et al., 2020).

Scene Segmentation. Later experiments require temporal segmentation of videos based on camera cuts. For this purpose, we run the dataset through the network of Souček and Lokoč (2020), which detects hard camera cuts. A scene, defined as the continuous shot between two cuts, lasts 2.2 seconds on average. This is similar to CMD, another movie dataset, whose scenes last 2.4 seconds on average. However, average scenes in ActivityNet (Caba Heilbron et al., 2015) and Kinetics-400 (Kay et al., 2017) last for 11.1 seconds and 30 seconds respectively. This shows camera cuts in movies are

much more frequent than the user-generated videos in ActivityNet and Kinetics.

5 Characteristics of SYMON Stories

5.1 Story Coverage

To facilitate story understanding, it is desirable that, despite their short lengths, the videos in SYMON provide sufficient coverage (Bain et al., 2020) over major plot points of the original movies or TV shows. In this experiment, we treat Wikipedia plot summaries (WikiPlots)¹ as ground truth and estimate the extent the stories in CMD, LSMDC, and SYMON cover the sentences in WikiPlots.

We use a three-step procedure for computing story coverage. First, we match movie summary in our dataset to their WikiPlots summaries by name. Second, we estimate if a sentence from the video narration is equivalent to a sentence in WikiPlots using the natural language inference (NLI) classifier from Nie et al. (2020). From two input sentences a and b , the NLI classifier predicts one of three possibilities: a entails b ; a contradicts b ; and neither is true. As entailment is asymmetric, we use the average probability for both directions (a entails b and b entails a) as the probability that a and b are equivalent. Finally, we find the best correspondence between two texts using Dynamic Time Warping (DTW) (Berndt and Clifford, 1994), which optimizes correspondence over entire sequences.

Briefly, DTW is a dynamic programming algorithm that seeks minimum-cost correspondence between two sequences, the WikiPlots sentence sequence A , and the narration sentence sequence B . We refer readers to the Appendix for a detailed description of the DTW algorithm. Using manually labeled sentence correspondences, we determine two model parameters, δ_A and δ_B , which denote the respective costs for skipping a sentence in sequences A and B .

We manually labeled the correspondence between around 500 sentences in CMD with Wikiplots stories, and did the same for SYMON. For LSMDC, we labeled around 1300 sentences because LSMDC texts are much longer. A second annotator labeled a small portion of data from each dataset to compute inter-rater reliability. The Cohen Kappa on SYMON, CMD and LSMDC are 0.86, 0.59, and 0.33 respectively. We believe the

¹<https://github.com/markriedl/WikiPlots>

	CMD	LSMDC	SYMOn
Story Coverage	10.8%	18.1%	37.9%

Table 2: Estimated story coverage with sentence entailment and Dynamic Time Warping.

low agreement on LSMDC is caused by the mismatch in the text lengths. Texts in LSMDC are longer than all other story texts, which led to difficulties in precisely locating the correspondence.

With a grid search, we find the optimal δ_A and δ_B as those that cause DTW to identify matched sentences the most accurately. The accuracy is defined as

$$\text{Accuracy} = \frac{1}{2} \left(\frac{T^{\text{wiki}}}{N^{\text{wiki}}} + \frac{T^{\text{text}}}{N^{\text{text}}} \right). \quad (1)$$

Here T^{wiki} and N^{wiki} are the number of correctly matched and the total number of WikiPlots sentences, respectively. T^{text} and N^{text} are the number of correctly matched and the total number of video narration sentences. We do not directly optimize story coverage because doing so results in incorrectly matched sentences that artificially inflate the story coverage measurement.

With the optimal δ_A and δ_B , we perform DTW again and calculate story coverage as the proportion of WikiPlots sentences matched with narration sentences,

$$\text{Coverage} = \frac{1}{K} \sum_i^K \frac{M_i}{N_i^{\text{wiki}}}, \quad (2)$$

where K is the number of WikiPlots movies appearing in the video dataset. In the i^{th} WikiPlots text, M_i denotes the number of matched sentences and N_i^{wiki} denotes the total number of sentences.

Table 2 shows the story coverage results. Of the three datasets, SYMON provides the highest coverage. LSMDC comes in second place, partially because it contains significantly longer descriptions for each movie than the other datasets.

5.2 Mental State Descriptions

A crucial component of story understanding is to develop theory of mind for the story characters, that is, to understand their mental states, such as emotions, motivations, and intentions (Bruner, 1986; Happé, 1994; Pelletier and Beatty, 2015). However, these concepts tend to be under-represented

	Emotion	Motivation	Intention
CMD	38.9	1.41	9.4
LSMDC	33.5	0.62	2.8
ActivityNet Captions	27.5	0.51	2.7
SYM _{ON} (Ours)	57.6	1.58	23.9

Table 3: Frequency of words related to emotion, motivation, and intention per one thousand words in the text corpora.

in the textual descriptions from commonly used video-language datasets.

In this experiment, we measure the frequency of words related to emotions, motivations, and intentions in the text associated with the videos. For emotional words, we adopt the WordNet-feelings dataset (Siddharthan et al., 2018), which includes 11387 emotion-related words identified by human experts. For motivation and intention words, we find 200 nearest neighbors of the words “motivation” and “intention” using 300-dimensional fast-Text embedding (Bojanowski et al., 2017) trained on Wikipedia and Common Crawl². We select 200 words as we find additional neighbors to be irrelevant to motivation and intention.

Table 3 reports word frequencies for every thousand words in four video-language datasets. We observe that SYMON employs mental-state words the most frequently and uses intention-related words 2.5 times as often as the next dataset, CMD. ActivityNet Captions (Krishna et al., 2017), containing matter-of-fact descriptions of actions in generic user-uploaded videos, uses the least of such words. LSMDC, which contains literal descriptions of movie clips, is ranked the third. CMD has a focus on the story content and is ranked the second. Overall, we find the ranking consistent with the nature of the datasets, as story text describes mental states more often than literal descriptions of generic videos. SYMON appears to be the most prototypical story dataset of the four.

6 Understanding Video-Text Divergence by Sequencing Videos

As discussed earlier, SYMON are characterized by large gaps between the textual and visual modalities due to the reporting bias, or the tendency to

²Acquired from <https://github.com/facebookresearch/fastText/blob/master/docs/crawl-vectors.md>.

avoid stating what can be easily observed from the video, and the prevalence of mental state descriptions, which are often not visible from the video. In this section, we report an experiment designed to estimate the extent of video-text correspondence.

Problem Definition. Similar to event/sentence ordering (Liu et al., 2018; Devlin et al., 2019), we predict the correct ordering of two consecutive video segments separated by a hard camera cut. The network predicts one of two classes: video segment 1 precedes segment 2 or vice versa. To create balanced classification, we randomly flip the ordering of the two video segments. We extract the text description that spans the same duration as the two video segments and expand the text to sentence boundaries.

We design two networks, one utilizing the unaltered textual description and the other solely relying on visual input. This setup allows us to estimate the amount of information provided by text. That is, if the text provides grounding to elements in both video segments, it should help the text-aware network predict the correct ordering.

Network Architecture. We adopt three pre-trained modules, the text encoder, the video encoder, and the cross-modality encoder from UniVL (Luo et al., 2020), which are pretrained on HowTo100M (Miech et al., 2019), and finetune the weights. The two video segments are encoded separately and their features are concatenated with the encoded text feature. After that, the two groups of features go through the cross encoder independently, yielding feature vectors f_1 and f_2 . With parameter w , the prediction is

$$P(\hat{y} = 1) = \sigma(w^\top f_1 - w^\top f_2). \quad (3)$$

where $\sigma(\cdot)$ is the sigmoid function and \hat{y} is the predicted class index. Figure 2 shows the overall network architecture.

As a baseline, we also create a network that relies on only the visual input, in which we replace the textual feature fed into the cross encoder with an all-zero vector. The rest of the network architecture remains the same.

Experimental Setup. To cover as much data as possible, we adopt a special dataset split, containing Set A of 2,444 videos, Set B of 2,289 videos, and a validation set of 500 videos. Each network is trained on Set A and tested on Set B, and then trained on Set B and tested on A. We report the

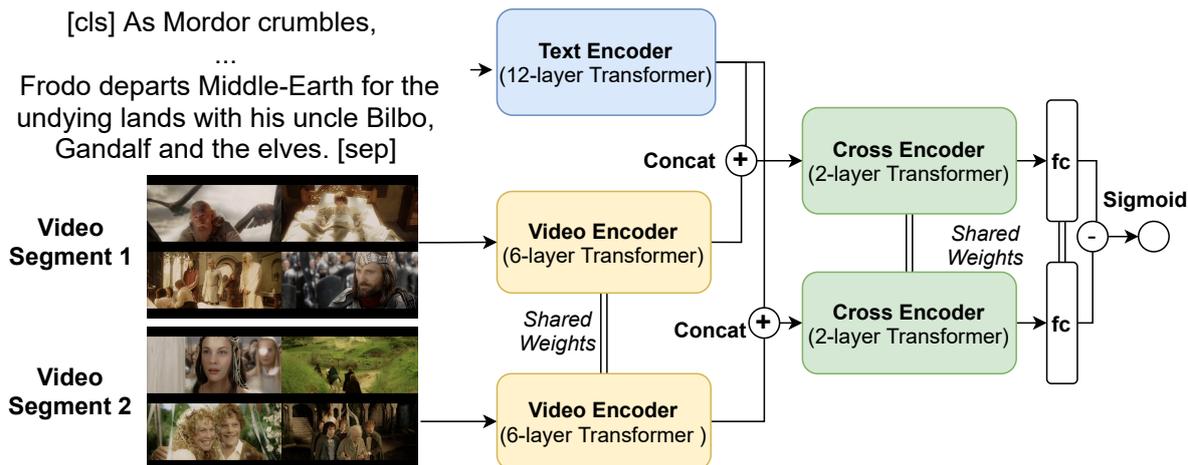


Figure 2: The network architecture for the temporal ordering task. The double vertical lines indicate weight sharing between modules.

average test accuracy. We tuned hyperparameters extensively on the validation set and select the training epoch with the highest validation accuracy. To avoid test data leak, we put all videos of the same movie or movie franchise to the same set. More settings can be found in the Appendix.

Results and Discussion. Table 4 lists the predictive accuracy. The network based solely on video has an accuracy of 63.4%. The incorporation of textual information improves prediction accuracy by 5.7%. Noting that chance level is at 50%, we find the performance to be mediocre. Since UniVL has been pretrained on HowTo100M and provides a good initialization, the results underscore the effects of the semantic gap between video and text.

Without text, 36.6% of data points cannot be correctly sequenced. Out of these, $5.7/36.6 = 15.6\%$ can be correctly classified with text. As the 36.6% are difficult data samples, we estimate the probability that (1) the text makes reference to both video segments *and* (2) the network correctly recognizes the references to be *at least* 15.6%.

Data Samples. In Figure 3, we present two data points, one from the 5% most helpful text cluster and one from the 5% least helpful text cluster. We observe that the helpful text mentions objects such as cauldron and book that appear in both video segments. As a result, both video segments can be grounded in the text, which provides ordering information. In comparison, the unhelpful text mentions rare object and action such as cat costume and jewelry robbery, which are difficult for the network to learn. Similarly, connecting the text “the mother

	Text + Video	Only Video
Accuracy	69.1%	63.4%

Table 4: Temporal order prediction accuracy of the text-aware and visual-only models.

refuses her son” and the discussion shown in video is not straightforward and would require identity tracking and event understanding.

Object and Action Analysis. We examine the match between text and video with contemporary technology on object detection and action recognition. First, for every data point, we compute the confidence of the ground-truth class from the two models. If the text-aware model has higher confidence than the visual-only model, we consider the text to be helpful. We rank the data points by the confidence difference between the two models, and take 5% data with the most helpful text and 5% data with the least helpful text.

Next, we run Faster-RCNN (Girshick, 2015) trained on Open Images V4 (Kuznetsova et al., 2020) to detect 600 object classes on video frames, and 3D-ResNet (Hara et al., 2018) trained on Kinetics-700 (Kay et al., 2017) to detect 700 action classes. After that, we match the identified objects and actions to the texts. The Appendix contains more details.

Table 5 shows that the most helpful texts contain relatively 18.8% more recognizable objects and 25.0% more actions than the most unhelpful texts. This suggests that textual references to ob-

Helpful Text



later on that night, elaine is in her apartment preparing a concoction of some sort, with ingredients being thrown into a small cauldron. she reads the ingredient list from an old apothecary book as she turns the page, we see that she is preparing for a love spell.

Unhelpful Text



a weirdo in a cat costume, walks in. he is actually the housekeeper's son, and comes there for shelter because he just robbed a jewelry store and escaped from the police. he wants the doctor to change his face to avoid being caught and sent to jail, but the mother refuses her son, believing that he's too crazy for that.

Figure 3: Examples from the most and least helpful text clusters. Bound boxes of the same color in text and video frame denote video-text correspondence. The black line denotes the boundary between the two video segments to be ordered.

	Objects Detected	Actions Recognized
Helpful Text	0.19	0.20
Unhelpful Text	0.16	0.16

Table 5: Number of words that match exactly the detected object names or action names per text description.

466 jects and actions in the video may have contributed
467 to the temporal ordering task. Noting that a text
468 description in this experiment contains 83 words
469 on average, the detected objects and actions ap-
470 pear rather scarce. We once again attribute this
471 observation to the reporting bias in the dataset.

7 Multimodal Retrieval

472 In this section, we establish baselines on the task of
473 video-text retrieval on SYMON and the YouTube
474 Movie Summary (YMS) (Dogan et al., 2018),
475 which serve as benchmarks for future research.
476

7.1 Network Architecture

477 We employ pretrained UniVL encoders without the
cross encoder. We encode the i^{th} text with the text
encoder, producing feature t_i , and encode the i^{th}
video segment with the video encoder, producing
feature v_i . The similarity between the two is simply
their dot product. With randomly sampled negative
text features $t_k, k \neq i$ and video features $v_k, k \neq i$,
we use the NCELoss (Gutmann and Hyvärinen,

2010):

$$L_{\text{NCE}} = \frac{1}{N} \sum_{i=1}^N -\mathbf{v}_i^\top \mathbf{t}_i + \log \left(\exp \mathbf{v}_i^\top \mathbf{t}_i + \sum_{k \neq i}^K \exp \mathbf{v}_i^\top \mathbf{t}_k + \sum_{k \neq i}^K \exp \mathbf{v}_k^\top \mathbf{t}_i \right) \quad (4)$$

478 where N is the total number of training samples
479 and K the number of negative samples.

7.2 Retrieval on SYMON

480 For the retrieval task, we create training, valida-
481 tion, and test sets with 4,191, 500 and 502 videos,
482 respectively. No movies or movie franchises ap-
483 pear in two sets simultaneously. The videos are
484 divided into non-overlapping clips, each consisting
485 of two scenes and having mean duration of 4.4 sec-
486 onds. YouTube videos often contain introduction
487 and channel information at the beginning and the
488 end, so we exclude 5% at each end of the videos.
489

490 In Table 6, we report recall at 1, 5, and 10 items
(R@1, R@5, and R@10), and Median Rank (MR).
491 As the video and the text are not exactly matched
492 by time, given a video clip, we consider the three
493 closest sentences as correct answers and vice versa.
494 As we expect, the UniVL network finetuned on
495 SYMON (UniVL-SYMON) outperforms the origi-
496 nal UniVL weights.
497

7.3 Transfer to YMS

498 Without in-domain finetuning, we directly test the
499 model trained on SYMON on the YMS dataset,
500 which contains 94 YouTube movie summary videos
501

Model	R@1	R@5	R@10	MR
<i>Text-to-video Retrieval</i>				
UniVL	0.11	0.39	0.63	4818
UniVL-SYMON	0.73	2.02	3.07	1785
<i>Video-to-text Retrieval</i>				
UniVL	0.03	0.11	0.19	5687
UniVL-SYMON	0.89	2.03	2.93	1843

Table 6: Retrieval performance on SYMON

with manual annotation of fine-grained video-text alignment. To prevent test data leak, we remove any summary videos for the 94 YMS movies from the training set used in this experiment.

Evaluation. In YMS, a text segment may correspond to multiple video clips, whereas a video clip may correspond to one or zero text segment. During inference, we align every video clip to the text segment with the highest similarity, as computed by the neural network. This creates the desired many-to-one alignment. If the highest similarity falls below a threshold, tuned on the validation set, the video clip is considered as not matching anything.

Following Dogan et al. (2018), we use clip accuracy (*i.e.* the temporal proportion of correctly aligned video segments), and sentence IoU (*i.e.* the intersection-over-union metric between aligned video durations and ground-truth durations) as evaluation metrics.

Baselines. Using the network described in §7.1, we compare the original UniVL weights, UniVL finetuned on SYMON data (UniVL-SYMON), as well as the supervised NeuMATCH network without the sequential context (*i.e.*, the minimum distance (MD) baseline from Dogan et al. (2018)). Note that UniVL-SYMON is trained with two video scenes as the basic unit for retrieval and NeuMATCH-MD uses more finely segmented units. As YMS contains fine-grained annotations, it is likely that this comparison puts our network at a disadvantage.

Test Data Split and Segmentation. For fair comparison with NeuMATCH-MD, we use the original test set of 15 videos and the original video segmentation. In addition, we also create a new split using 70% of the entire YMS as the test set and 30% as the validation set. In this new setting, the

	Clip Acc.	Sent. IoU
<i>Original Data Split and Segmentation</i>		
UniVL	3.7	1.5
NeuMATCH-MD (Supervised)	4.0	2.4
UniVL-SYMON	5.4	2.6
<i>New Data Split and Segmentation</i>		
UniVL	4.3	2.1
UniVL-SYMON	6.2	2.4

Table 7: Zero-shot alignment performance on YMS.

videos are segmented into scenes as detected in our preprocessing (§4).

Results. Table 7 shows the results. Despite the difference in segmentation and the weak supervision from SYMON, UniVL-SYMON outperforms the supervised NeuMATCH-MD baseline. This shows that UniVL-SYMON learns a superior cross-modality distance metric, demonstrating the utility of the large-scale SYMON dataset. UniVL-SYMON also outperforms the original UniVL by 1.7% / 1.1% in the original setting and 1.9% / 0.3% in the new setting. Considering UniVL was trained on the gigantic HowTo100M dataset, we attribute the improvement to the similarity between SYMON and YMS, which highlights the effectiveness of SYMON in the domain of story video understanding.

8 Conclusion

In this work, we collect and process a story understanding SYMON. We compare SYMON with existing video-language datasets and quantitatively analyze the story coverage, the amount of mental-state descriptions, and the semantic divergence between video and text. Furthermore, we establish multimodal retrieval baselines for SYMON and a zero-shot alignment baseline on YMS to demonstrate the effectiveness of SYMON in story understanding. We believe SYMON will serve as a new challenge for the research community and inspire new advances of multimodal machine learning.

9 Potential Ethical Impact

In this paper, we collect user-uploaded videos from YouTube, which are summaries of mostly western movies and TV shows in the English language. We recognize that movies and TV shows are fictional in nature, and often prioritize dramatic events over

faithful representation of real-life scenarios. In addition, the videos may reflect particular bias of the creators of the movie and TV shows or the creators of the summary videos, as well as bias from particular cultures or the time periods of production.

For these reasons, we urge researchers to take caution when attempting to learn social norms from such videos. For example, events of bank robberies may be over-represented in these videos, and a machine learning model may inadvertently infer that robbing a bank is part of the social norm. In addition, the model may incorrectly learn from disproportional association of certain groups of people with certain social status, occupations, and other cultural constructs.

We further note that most relations between events are probabilistic and neither necessary nor sufficient. For example, though it is common for someone with a medical emergency to call for an ambulance, it does not always happen. We suggest researchers to similarly qualify any learned relations. The dataset is intended for fundamental research and not real-world deployment.

References

Tanvirul Alam, Akib Khan, and Firoj Alam. 2020. Punctuation restoration using transformer models for resource-rich and-poor languages. In *Proceedings of the Sixth Workshop on Noisy User-generated Text (W-NUT 2020)*, pages 132–142.

Youngmin Baek, Bado Lee, Dongyoon Han, Sangdoon Yun, and Hwalsuk Lee. 2019. Character region awareness for text detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9365–9374.

Max Bain, Arsha Nagrani, Andrew Brown, and Andrew Zisserman. 2020. Condensed movies: Story based retrieval with contextual embeddings. In *Proceedings of the Asian Conference on Computer Vision*.

Donald J. Berndt and James Clifford. 1994. Using dynamic time warping to find patterns in time series. In *KDD workshop*.

Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Scott Wen tau Yih, and Yejin Choi. 2020. Abductive commonsense reasoning. In *ICLR*.

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. "O'Reilly Media, Inc."

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.

Faeze Brahman, Meng Huang, Oyvind Tafjord, Chao Zhao, Mrinmaya Sachan, and Snigdha Chaturvedi. 2021. "let your characters tell their story": A dataset for character-centric narrative understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Jerome Bruner. 1986. *Actual Minds, Possible Worlds*. Harvard University Press.

Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Nieves. 2015. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970.

Tommaso Caselli, Eduard Hovy, Martha Palmer, and Piek Vossen. 2021. *Computational Analysis of Storylines: Making Sense of Events*. Cambridge University Press.

Nathanael Chambers, Taylor Cassidy, Bill McDowell, and Steven Bethard. 2014. [Dense event ordering with a multi-pass architecture](#). *Transactions of the Association for Computational Linguistics*, 2:273–284.

Nathanael Chambers and Dan Jurafsky. 2008. Unsupervised learning of narrative event chains. In *Proceedings of ACL-08: HLT*.

Snigdha Chaturvedi, Shashank Srivastava, Hal Daume III au2, and Chris Dyer. 2016. Modeling dynamic relationships between characters in literary novels. In *AAAI*.

Mingda Chen, Zewei Chu, Sam Wiseman, and Kevin Gimpel. 2021. Summscreen: A dataset for abstractive screenplay summarization. *arXiv Preprint 2104.07091*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Pelin Dogan, Boyang Li, Leonid Sigal, and Markus H. Gross. 2018. [LSTM stack-based neural multi-sequence alignment technique \(neumatch\)](#). *CoRR*, abs/1803.00057.

David Elson, Nicholas Dames, and Kathleen McKeown. 2010. Extracting social networks from literary fiction. In *Proceedings of the 48th annual meeting of the*

790	Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac,	Matthew Richardson, Christopher J.C. Burges, and Erin	845
791	Makarand Tapaswi, Ivan Laptev, and Josef Sivic.	Renshaw. 2013. MCTest: A challenge dataset for the	846
792	2019. Howto100m: Learning a text-video embed-	open-domain machine comprehension of text . In <i>Pro-</i>	847
793	ding by watching hundred million narrated video	<i>ceedings of the 2013 Conference on Empirical Meth-</i>	848
794	clips. In <i>Proceedings of the IEEE/CVF International</i>	<i>ods in Natural Language Processing</i> , pages 193–203,	849
795	<i>Conference on Computer Vision</i> , pages 2630–2640.	Seattle, Washington, USA. Association for Computa-	850
		tional Linguistics.	851
796	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal,	Melissa Roemmele, Cosmin Adrian Bejan, and And-	852
797	Jason Weston, and Douwe Kiela. 2020. Adversarial	rew S. Gordon. 2011. Choice of plausible alter-	853
798	NLI: A new benchmark for natural language under-	natives: An evaluation of commonsense causal rea-	854
799	standing . In <i>Proceedings of the 58th Annual Meet-</i>	soning. In <i>AAAI Spring Symposium on Logical For-</i>	855
800	<i>ing of the Association for Computational Linguistics</i> ,	<i>malizaciones of Commonsense Reasoning</i> .	856
801	pages 4885–4901, Online. Association for Computa-		
802	tional Linguistics.		
803	Qiang Ning, Hao Wu, Rujun Han, Nanyun Peng, Matt	Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and	857
804	Gardner, and Dan Roth. 2020. Torque: A reading	Bernt Schiele. 2015. A dataset for movie description.	858
805	comprehension dataset of temporal ordering ques-	In <i>Proceedings of the IEEE conference on computer</i>	859
806	tions. In <i>The 2020 Conference on Empirical Methods</i>	<i>vision and pattern recognition</i> , pages 3202–3212.	860
807	<i>in Natural Language Processing (EMNLP)</i> , pages		
808	1158–1172. Association for Computational Linguis-	Anna Rohrbach, Atousa Torabi, Marcus Rohrbach,	861
809	tics.	Niket Tandon, Christopher Pal, Hugo Larochelle,	862
		Aaron Courville, and Bernt Schiele. 2017. Movie	863
		description. <i>International Journal of Computer Vi-</i>	864
		<i>sion</i> , 123(1):94–120.	865
810	Tim O’Gorman, Kristin Wright-Bettner, and Martha	Arka Sadhu, Tanmay Gupta, Mark Yatskar, Ram Neva-	866
811	Palmer. 2016. Richer event description: Integrating	tia, and Aniruddha Kembhavi. 2021. Visual seman-	867
812	event coreference with temporal, causal and bridging	tic role labeling for video understanding . <i>CoRR</i> ,	868
813	annotation . In <i>Proceedings of the 2nd Workshop on</i>	abs/2104.00990.	869
814	<i>Computing News Storylines (CNS 2016)</i> , pages 47–		
815	56, Austin, Texas. Association for Computational		
816	Linguistics.		
817	Jessica Ouyang and Kathleen McKeown. 2015. Mod-	Maarten Sap, Ronan LeBras, Emily Allaway, Chan-	870
818	eling reportable events as turning points in narrative .	dra Bhagavatula, Nicholas Lourie, Hannah Rashkin,	871
819	In <i>Proceedings of the 2015 Conference on Empiri-</i>	Brendan Roof, Noah A. Smith, and Yejin Choi. 2019.	872
820	<i>cal Methods in Natural Language Processing</i> , pages	Atomic: An atlas of machine commonsense for if-	873
821	2149–2158, Lisbon, Portugal. Association for Com-	then reasoning. In <i>The AAAI Conference on Artificial</i>	874
822	putational Linguistics.	<i>Intelligence</i> .	875
823	Pinelopi Papalampidi, Frank Keller, Lea Frermann, and	Advaith Siddharthan, Nicolas Cherbuin, Paul J Eslinger,	876
824	Mirella Lapata. 2020. Screenplay summarization us-	Kasia Kozłowska, Nora A Murphy, and Leroy Lowe.	877
825	ing latent narrative structure . In <i>Proceedings of the</i>	2018. Wordnet-feelings: a linguistic categorisation	878
826	<i>58th Annual Meeting of the Association for Compu-</i>	of human feelings. <i>arXiv preprint arXiv:1811.02435</i> .	879
827	<i>tational Linguistics</i> , pages 1920–1933, Online. Asso-		
828	ciation for Computational Linguistics.	Gunnar A Sigurdsson, Gül Varol, Xiaolong Wang, Ali	880
829	Janette Pelletier and Ruth Beatty. 2015. Children’s	Farhadi, Ivan Laptev, and Abhinav Gupta. 2016. Hol-	881
830	understanding of aesop’s fables: relations to reading	lywood in homes: Crowdsourcing data collection for	882
831	comprehension and theory of mind. <i>Frontiers in</i>	activity understanding. In <i>European Conference on</i>	883
832	<i>Psychology</i> , 6:1448.	<i>Computer Vision</i> , pages 510–526. Springer.	884
833	Karl Pichotta and Raymond J. Mooney. 2016. Using	Tomáš Souček and Jakub Lokoč. 2020. Transnet v2: An	885
834	sentence-level LSTM language models for script in-	effective deep network architecture for fast shot tran-	886
835	ference . In <i>Proceedings of the 54th Annual Meeting</i>	sition detection. <i>arXiv preprint arXiv:2008.04838</i> .	887
836	<i>of the Association for Computational Linguistics (Vol-</i>		
837	<i>ume 1: Long Papers)</i> , pages 279–289, Berlin, Ger-	Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen,	888
838	many. Association for Computational Linguistics.	Antonio Torralba, Raquel Urtasun, and Sanja Fi-	889
839	Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A.	dler. 2016. MovieQA: Understanding Stories in	890
840	Smith, and Yejin Choi. 2018. Event2Mind: Com-	Movies through Question-Answering. In <i>IEEE Con-</i>	891
841	monsense inference on events, intents, and reactions.	<i>ference on Computer Vision and Pattern Recognition</i>	892
842	In <i>Proceedings of the 56th Annual Meeting of the</i>	<i>(CVPR)</i> .	893
843	<i>Association for Computational Linguistics (Volume</i>	Atousa Torabi, Christopher Pal, Hugo Larochelle, and	894
844	<i>1: Long Papers)</i> .	Aaron Courville. 2015. Using descriptive video ser-	895
		vices to create a large data source for video annota-	896
		tion research. <i>arXiv preprint arXiv:1503.01070</i> .	897

Naushad UzZaman, Hector Llorens, Leon Derczynski, James Allen, Marc Verhagen, and James Pustejovsky. 2013. [SemEval-2013 task 1: TempEval-3: Evaluating time expressions, events, and temporal relations](#). In *Second Joint Conference on Lexical and Computational Semantics (*SEM): the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pages 1–9, Atlanta, Georgia, USA. Association for Computational Linguistics.

Siddharth Vashishtha, Adam Poliak, Yash Kumar Lal, Benjamin Van Durme, and Aaron Steven White. 2020. Temporal reasoning in natural language inference. In *Findings of the Association for Computational Linguistics: EMNLP 2020*.

Siddharth Vashishtha, Benjamin Van Durme, and Aaron Steven White. 2019. [Fine-grained temporal relation extraction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2906–2919, Florence, Italy. Association for Computational Linguistics.

Paul Vicol, Makarand Tapaswi, Lluís Castrejon, and Sanja Fidler. 2018. Moviegraphs: Towards understanding human-centric situations from videos. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan Roth. 2021. Joint constrained learning for event-event relation extraction. In *EMNLP*.

Chao-Yuan Wu and Philipp Krahenbuhl. 2021. Towards long-form video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1884–1894.

Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. 2018. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *Proceedings of the European conference on computer vision (ECCV)*, pages 305–321.

Zhengzhe Yang and Jinho D. Choi. 2019. [FriendsQA: Open-domain question answering on TV show transcripts](#). In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 188–197, Stockholm, Sweden. Association for Computational Linguistics.

Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot, Ashish Sabharwal, and Dan Roth. 2021. Temporal reasoning on implicit events from distant supervision. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

A Story Coverage

Dynamic Time Warping We present the DTW problem formulation: given the WikiPlots sequence of sentences $A = (a_1, \dots, a_N)$ and the video narration sentences $B = (b_1, \dots, b_M)$, we

seek the best set of correspondence $\{(a_i, b_{g(i)})\}_{i=1}^N$, where the function $g(i) \in \{\epsilon, 1, \dots, M\}$ returns the index in sequence B that matches sentence a_i in A . Setting $g(i) = \epsilon$ indicates that a_i is not matched with any sentence in B .

The DTW algorithm can be understood as finding the shortest path in a graph, where each node (i, j) in the graph represents matching sentence a_i and sentence b_j . The graph contains dummy nodes $(0, 0)$ and $(N + 1, M + 1)$. From node (i, j) , we can transit to node $(i + 1, j + 1)$, which would match a_{i+1} with b_{j+1} and incur cost $c(i + 1, j + 1)$.

$$c(i + 1, j + 1) = 1 - P(a_{i+1} \Leftrightarrow b_{j+1}). \quad (5)$$

Here $P(a_{i+1} \Leftrightarrow b_{j+1})$ denotes the probability that sentences a_{i+1} and b_{j+1} are equivalent, as determined by the Natural Language Inference classifier.

Similarly, we can transit from (i, j) to $(i + 1, j)$, which would match a_{i+1} with b_j and incur cost $c(i + 1, j)$. The transition from (i, j) to $(i, j + 1)$ is symmetric. Additionally, we can transit from (i, j) to $(i, j + 1, \epsilon)$, which prevents b_{j+1} from matching anything. From $(i, j + 1, \epsilon)$, we may transit to $(i, j + 2, \epsilon)$, $(i, j + 2)$, or $(i + 1, j + 2)$. The costs of ignoring a sentence in A and B are δ_A and δ_B respectively. With this setup, the best correspondence can be found as the path from $(0, 0)$ to $(N + 1, M + 1)$ with minimum cost. We find optimal δ_A and δ_B using manually labeled sentence correspondence.

Annotation instructions Fig. 4 shows the instructions we give to our annotator. Here column A is the WikiPlot summary and column B is the summary from SYMON or CMD or LSMDC.

B Video Temporal Ordering

Hyperparameters. We sample each video segment at 16 frames per second (FPS) and extract features with S3D (Xie et al., 2018) pretrained on HowTo100M. Each video video segment last exactly 8 seconds. We extract S3D features every second (i.e. from 16 frames), yielding 8 1024-dimensional video features for each video segment. For video features extraction we use frame size of 112×112 .

We extract the text between the start of the first video segment and the end of the second video segment. To ensure completeness, the text is extended to the nearest sentence boundaries. The maximal

Instructions

Columns **A** and **B** are different narratives of the same story. For each sentence in column **A** try to find an equivalent sentence in column **B** and put its index in the brackets. If there's not equivalent sentence to be found, leave the bracket empty.

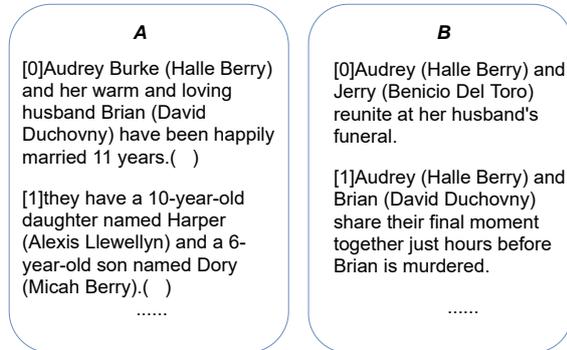


Figure 4: Annotation Instruction

number of text tokens is 128. For longer texts, we remove extra tokens from the start and end of the text. For shorter texts, we add zero padding to the end.

The text encoder, video encoder, and cross encoder consist of 12, 6, 2 Transformer layers, respectively. The models are trained for 30 epochs with learning rate warm-up in the first 6 epochs. Hyperparameters are tuned on the validation set. The text-aware model is trained with a batch size of 128 and learn rate of $5e - 6$ and the visual-only model is trained with a batch size of 256 and initial learning rate of $1e - 5$. We apply cosine learning rate decay and the Adam optimizer to all models.

The model contains 217,185,539 parameters and is trained for 2.7 hours on 4 Nvidia 3090 GPUs. The results reported in the main paper are on a single run.

Calculating overlap between text description and object/action class

We first tokenize the text description and use part-of-speech tagging to identify nouns and verbs in the text description (Bird et al., 2009). For matching with object and action detection, we retain the nouns and verb from text description, respectively. We also lemmatize the retained words to remove variations and remove common nouns and verbs (“men”, “women”, “person”, and “clothing” for nouns and “is”, “go”, “to”, “get”, “have”, “look”, “walk”, “play”, and “take” for verbs). For object detection we retain the top 10 class predictions for each clip. For action detection we divide the clip into scenes and retain the

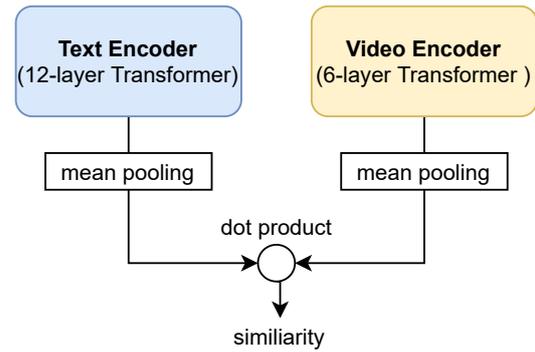


Figure 5: Retrieval Model

top 3 class prediction for each scene. Finally, we calculate the number of time the retained nouns or verbs appear in the detected object or action class names.

C Retrieval

Hyperparameters. For video feature extraction we sample the video at 16 FPS and use S3D pretrained on HowTo100m to extract one 1024 dimensional feature every 16 frames. The frame size is 112×112 . For each clip we extract 4 features, if the clip is shorter than 4 seconds zero padding is add and if the clip is longer than 4 seconds we only use the first 4 second. Likewise, we take 64 text tokens for each clip. Text is extracted from between the start and end of the video clip and extended to the nearest sentence boundaries. The video and text encoders consist of 12 and 6 transformer layers respectively, and are initialized from UniVL pretrained on HowTo100m. The outputs are then averaged into two 768 dimensional embeddings for video and text. The similarity between a video, text pair is calculated as the dot product of the video and text embeddings. The model is finetuned on SYMON with an initial learning rate of $5e - 5$ and cosine learning rate decay. We use a batchsize of 1024 and train for 20 epoches, the first epoch is warm up. SGD with momentum of 0.9 is used for optimization and a weight decay term of 0.5 is added for regularization.

The model contains 153,784,064 parameters and is trained for 4 hours on 4 Nvidia 3090 GPUs. The results reported in the main paper are on a single run.

D Implementation and Licensing Details

For the subtitle masking in §4 we used EasyOCR (Baek et al., 2019) for image character recogni-

tion. The model we use is the english_g2 model from <https://www.jaided.ai/easyocr/modelhub/>. EasyOCR (Baek et al., 2019) is licensed under the Apache License, Version 2.0.

For the punctuation restoration in §4 we used the network from Alam et al. (2020). The model we use is given here <https://drive.google.com/file/d/17BPcnHVhpQlsOTC8LEayIFFJ7WkL00cr/view>. The network and model are released under the MIT license.

For scene segmentation in §4 we used TransNet-V2 (Souček and Lokoč, 2020) to identify scene boundaries, the network weight are from think link <https://github.com/soCzech/TransNetV2/tree/master/inference/transnetv2-weights>. For every frame, the network predict the probability of a scene change occurring immediately after the frame, if the probability is larger than a threshold of 50%, we deem a scene change had occurred. TransNet-v2 (Souček and Lokoč, 2020) is released under the MIT license.

For entailment prediction in §5.1 we use AdversarialNLI (Nie et al., 2020), specifically the 'roberta-large-snli_mnli_fever_anli_R1_R2_R3-nli' model. AdversarialNLI (Nie et al., 2020) is released under the MIT licence. For this section we use WikiPlot summaries from <https://github.com/markriedl/WikiPlots> as ground truth movie summaries. The release does not include a license. Additionally, we compare our dataset to the CMD (Bain et al., 2020) dataset and LSMDC (Rohrbach et al., 2017) dataset, both are released under the Creative Commons Attribution 4.0 International License.

For the mental-state description experiment in §5.2, we collect emotion related words from WordNet-feelings (Siddharthan et al., 2018) dataset. The release does not include a license. We collect intention, motivation related words from the top 200 nearest neighbors on Fast-text (Girshick, 2015) word embedding, which is acquired from <https://github.com/facebookresearch/fastText/blob/master/docs/crawl-vectors.md>. fast-Text (Girshick, 2015) is released under the MIT licence. The word embedding is release under the Creative Commons Attribution-Share-Alike License 3.0. The ActivityNet dataset (Krishna et al., 2017) is licensed under the MIT license.

For the video sequencing experiment in §6, we use UniVL (Luo et al., 2020) pretrained on HowTo100m (Miech et al., 2019). The model weights are initialized from <https://github.com/microsoft/UniVL/releases/download/v0/univl.pretrained.bin>. UniVL (Luo et al., 2020) and HowTo100m (Miech et al., 2019) are licensed under MIT and Apache License 2.0 respectively.

For object recognition in §6, we use Faster-RCNN (Girshick, 2015) trained on Open Images V4 (Kuznetsova et al., 2020) to detect objects from video frames, and 3D-ResNet (Hara et al., 2018) trained on Kinetics-700 (Kay et al., 2017) to detect actions. Faster-RCNN (Girshick, 2015) and 3D-ResNet (Hara et al., 2018) are licensed under the MIT license. Open Images V4 is released under Apache License 2.0. Kinetics-700 is licensed under the Creative Commons Attribution 4.0 International License. The text descriptions are processed with the nltk package (Bird et al., 2009), licensed under Apache License 2.0.

For the multimodal retrieval task in §7, we use the YMS dataset (Dogan et al., 2018) from <https://github.com/RubbyJ/Data-efficient-Alignment>.