

## Appendix

Errata: In Figure 1, Event 2 Arg2 should be “man with trident” instead of “main with trident”.

Appendix provides details on:

1. A Brief Summary of Semantic Roles, and their usage in our paper.
2. Details on Dataset Curation and Annotation Interface
3. Additional Dataset Statistics
4. Additional Implementation Details
5. Details on Lea-Soft along with Tables with All Metrics
6. DataSheet [18] for VidSitu
7. Qualitative Analysis of Data (this is attached as a video file in the zip folder).

### A. Semantic Roles: A Brief Summary

Semantic Role Labeling attempts to abstract out at a high-level who does what to whom [66]. It is a popular natural language task which attempts at obtaining such structured outputs from natural language descriptions. As such there are multiple sources to obtain semantic roles such as FrameNet [4], PropBank [54] and VerbNet [7]. Prior work on situation recognition in images (ImSitu) [83] have curated list of verbs (situations) from FrameNet, and action recognition dataset (Moments in Time) [51] have curated action vocabulary from VerbNet. However, we qualitatively found both vocabulary to be insufficient to represent actions, and thus chose PropBank which contained action-oriented verbs. As such, PropBank has been used for video object grounding [61] but not in the context of collecting semantic roles from visual data.

PropBank contains a set of numbered semantic roles for each verb ranging from Arg0 to Arg4. Each numbered argument has a specific definition for a particular verb but some themes are similar across verbs (adapted from PropBank annotation guidelines [6]<sup>3</sup>). For the verb “throw”:

- Arg0: Agent – object performing the action. For *e.g.* “person”
- Arg1: Patient – object on which action is performed. For *e.g.* “ball”
- Arg2: Instrument, Benefactive, Attribute. For *e.g.* “towards a basket”
- Arg3: Starting Point
- Arg4: Ending Point
- ArgM: Modifier – location (LOC), manner(MNR), direction (DIR), Purpose (PRP), Goal (GOL), Temporal (TMP), Adverb (ADV)

<sup>3</sup>[http://clear.colorado.edu/compsem/documents/propbank\\_guidelines.pdf](http://clear.colorado.edu/compsem/documents/propbank_guidelines.pdf)

In general, we noticed that Arg3 and Arg4 were exceedingly rare for visual verbs, thus we restrict our attention to Arg0, Arg1, Arg2 for numbered arguments. For modifier arguments, we found Location (LOC) to be universally valid for all video segments. Thus, for those verbs where LOC doesn’t apply usually, we additionally add a semantic role “Scene” which refers to “where” the event takes place (such as “living room”, “near a lake”). Other arguments were chosen based on their appearance in MPIID dataset, and we most commonly used Manner (which suggests “how” the action takes place) and Direction (details in the Section B). For rest of the paper, we use ALoc, ADir, AMnr, and AScn to denote location, direction, manner and scene arguments respectively.

### B. Dataset Collection

In this section we describe details on dataset collection including curation of verbs and arguments, followed by details on annotation interface, quality control and reward structure.

#### B.1. Dataset Curation

We provide more details on Dataset Curation which were omitted from Section 4.1 of the main paper.

**Video Source Selection.** As suggested in the Section 4.1 we aimed at a domain with two criterion: the videos should be by themselves cover diverse situations (“climb” verb should not just be associated with rocks or mountains, but also things like top of a car), and that the each video should contain complex situation (the video shouldn’t depict someone doing the same task over extended period of time, which would lower chances of finding meaningful event relations and be repetitive in verbs and arguments over the entire video).

After a brief qualitative analysis, we found instruction domain videos (HowTo100M [50], YouCookII [90], COIN [69]) to have very fine-grained actions with less diversity and less complexity within small segments, open domain sources (ActivityNet [24], Moments in Time [51], Kinetics [31], HACS[87]) to be somewhat diverse but low complexity within a small segment. This led us to Movie domain which span multiple genres leading to appreciable diversity as well as complexity.

We converged on using MovieClips [3] rather than other movie sources such as MPII [60], since MovieClips already provide one-stage of filtering to provide interesting videos. While using the same movies as used in AVA[21] was an option, we found that the video retention was quite low (around 20% of the movie are removed from youtube), and the movie contained long contiguous segments with low complexity. We also note some other datasets like MovieNet [28], Movie Synopsis Dataset [80], Movie Graphs [72] do not provide movie videos and cannot be

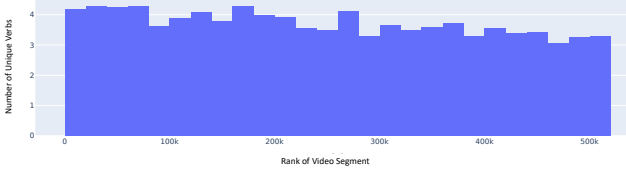


Figure 1: Bar graph showing number of unique verbs with respect to the rank of the video segment as computed via our heuristic based on predicted labels from SlowFast Network [16] trained on AVA[21].

used for collecting annotations. One demerit of using movie domain is that the verb distributions are skewed towards actions like “talk”, “walk”, “stare”. Despite this we find the videos to be reasonably complex.

**Video Selection.** MovieClips spans a total of 1k Hours which is far beyond what can be reasonably annotated. To best utilize available annotation budget, we are primarily interested in identifying video segments depicting complex situations with a high precision while avoiding visually uneventful segments common in movies such as those simply involving actors engaged in dialogue.

To avoid such segments, we use the following heuristic: a video with more atomic actions per person is likely to be more eventful. So, we divide all movieclips into 10 second videos with a stride of 5 seconds, obtain human bounding boxes from the MaskRCNN [23] object detector trained on the MSCOCO [45] dataset, predict atomic actions for each detected person using the SlowFast [16] activity recognition model trained on the AVA [21] dataset, and rank all videos by the average number of unique atomic actions per person in the video. In particular, we discard labels such as “talk”, “listen”, “stand” and “sit” as these atomic actions didn’t correlate with complexity of situations. Since “action” sequences like “fight scenes” are favored by our ranking measure, we use simple heuristic of removing “martial arts” actions to avoid oversampling such scenes and improve diversity of situations represented in the selected videos.

To confirm the usefulness of the proposed heuristic, we conduct an experiment where we annotate 1k videos chosen uniformly sampled across the entire dataset (as shown in Figure 1). Reducing number of unique verbs shows the effectiveness of our heuristic and suggests at least 80K videos segments (which translates to 27K non-overlapping video segments) can be richly annotated.

For final video selection, we randomly choose set of videos from the top-K ranks, such that the newly chosen videos don’t overlap with already chosen videos, and that no more than 3 videos are uploaded from the same Youtube video within a particular batch.

**Curating Verb Senses.** To curate verb senses, we follow a two-step process: from the initial list of  $\sim 6k$  verb senses

in PropBank [54], first we manually filter verb senses which share the same lemmatized verb (as previously stated “go” has 23 verb senses) to retain only “visual” verb senses (for instance we remove the verb sense of “run” which refers to running a business). We keep all 3.7K verbs with a single verb sense and of the remaining 2364 verbs-senses (shared across 809 verbs) we retain 629 verb senses (shared across 561 verbs). Second, to further restrict the set of verbs to those useful for describing movies, we discard verbs that do not appear at all in the MPII-Movie Description (MP2D) dataset [60]. To extract verbs from the descriptions we use a semantic-role parser [62]. This results in a final set of 2154 verb-senses.

**Curating Argument Roles.** Once we have curated the verb-senses from PropBank, we aim to delegate a set of argument roles for each verb-sense which would be filled based on the video. While PropBank provides numbered arguments for each verb-sense there are two issues with directly using them: first, some arguments are less relevant for visual scenes (for instance Arg1 (utterance) for “talk” is not visual), second, auxiliary arguments like direction and manner are not provided (for instance direction and manner for “look” are important to describe a scene). To address this issue, we re-use the MP2D dataset to inform us what arguments are used with the verbs. For each verb, we choose set of 5 most frequently used argument role-set and use their union. We also remove roles such as TMP (usually referring to words like “now”, “then”) since temporal context is implicit in our annotation structure. We also removed roles like ADV (adverb) which were too infrequent. Finally, we use the following modifier roles: “Manner”, “Location”, “Direction”, “Purpose”, “Goal”, but note that “purpose” and “goal” were restricted to a small number of verbs and hence not considered for evaluation.

We further added the modifier role “Scene” which describes “where” the event takes place, and only applies to verbs which don’t have “Location”. For instance, “stand” has the argument role “location” which refers to “where” the person is standing and doesn’t have “Scene”, whereas “run” doesn’t contain “location” and hence contains “Scene”. In general, “Scene” refers to the “place” of the event such as “in an alleyway” or “near a beach”.

**Event Relations.** We started with the set of three event relations namely: no relation (Events A and B are unrelated), causality (Event B is Caused By Event A *i.e.* B happens directly as a result of A) and contingency based (Event B is Enabled By Event A *i.e.* A doesn’t directly cause B but B couldn’t have happened without A happening first) on prior work in cross-document event relations [26]. However, we found adding an additional case of “Reaction To” for causality helpful to distinguish between event relations. For instance, in the case “X punches Y” followed by “Y falls down” would be definitely “B is Caused By A”, how-

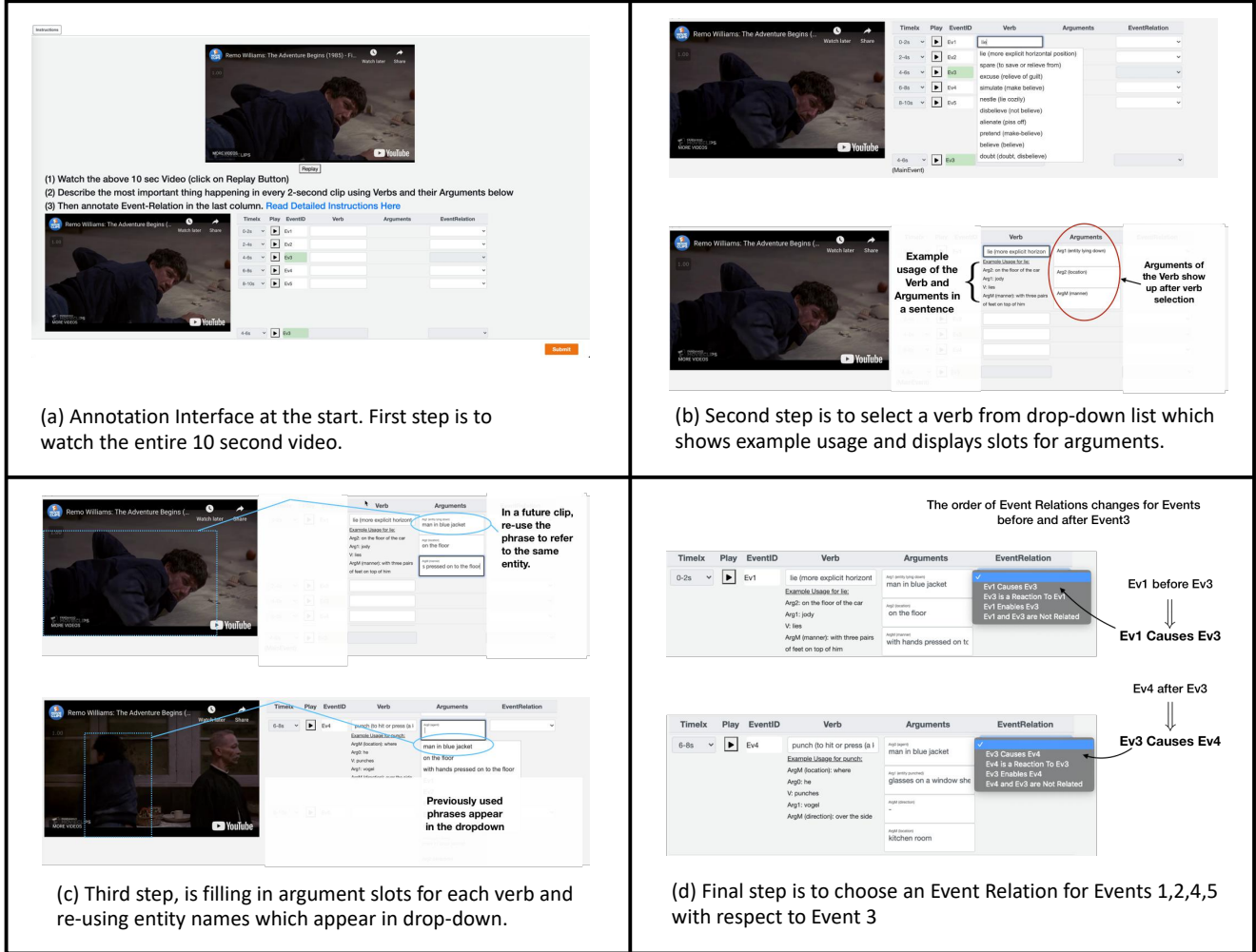


Figure 2: Illustration of our annotation interface. (a) depicts the initial screen an annotator sees. In the first step, one needs to watch the entire 10 second video. (b) depicts the second step of choosing a verb from a drop-down which contains verb senses obtained from PropBank. After selecting a verb, an example usage is shown along with corresponding argument roles which need to be filled. (c) depicts filling the argument slots for each verb which can be phrases of arbitrary length. Each filled in phrase can be re-used in a subsequent slot, to enforce co-reference of the entities. (d) shows the final step of choosing event relations once all the arguments for all events are filled. The event relations should be classified based on causality and contingency for Events 1,2,4,5 with respect to Event 3.

ever for the case “X punches Y” followed by “Y crouches” it is unclear if “B is Caused By A” since Y makes a voluntary decision to crouch. As a result, we call this relation “B is a Reaction To A”.

## B.2. Annotation pipeline

With videos, the list of verb-sense and their roles curated, we are now ready to crowd-source annotations on Amazon Mechanical Turk (AMT).

**Annotation Interface.** Figure 2 shows screenshots depicting our annotation interface. For annotating a given 10 second video, the assigned worker is instructed to first

watch the entire 10-second video (Figure 2 (a)). Then for every 2 second interval, the annotator selects a verb corresponding to the most salient event from our curated list of verb-senses using a search-able drop-down menu. Once the verb is chosen, slots for the corresponding roles are displayed along with an example usage (Figure 2 (b)). The worker fills in the values for each role using free-form text (typically a short phrase). When referring to an entity, we instruct the worker to use phrases that uniquely identify the entity in the full 10 second video. Furthermore, these phrases can be reused in filling semantic-roles in other events within the video, which provides the co-reference in-

formation about the entities *i.e.* co-referenced entities are maintained via exact-string match (Figure 2 (c)). Once all verbs and their roles are annotated, we ask the worker to label the relation of Events 1, 2, 4, and 5 with respect to Event 3 (Figure 2 (d)). Note that the order of causality and contingency is different for Events 4,5 compared to Events 1,2 respecting the temporal order.

**Partitioning into 2-second clips:** We emphasize that splitting the video into 2-second intervals is strictly a design choice motivated by reduction in annotation cost and consistent quality of annotations. In an early version of the data collection, we asked annotators to provide “start” and “end” points for events and allowed overlaps (consistent with other datasets such as ActivityNet Captions[35]). A close analysis showed that the noise in annotations was tremendous, took significantly longer (roughly 3x) and would lead to a much smaller and lower quality dataset given a budget. We thus simplified the task via 2-sec interval annotations and saw large improvements in consensus and speed.

Clearly, using such a scheme leads to imprecise temporal boundaries for the events. Furthermore, it doesn’t allow annotating hierarchical actions. However, we argue that *the downsides of this design choice are reasonably mitigated* since: (a) Longer duration events get annotated via a *repeat* of the same verb across consecutive clips (we see many occurrences in our dataset) & (b) In the presence of multiple verbs in a clip, the most *salient* one gets annotated.

The 2s duration was chosen after an analysis of  $\sim 50$  videos showed that events typically spanned more than 1s but clips longer than 2s often contained multiple interesting events that we would not want to discard. Finally, we note that 2-second duration choice may not be suitable for vastly different domains (e.g. fewer actions and more talking) where 2s may be too dense, and relaxing this to longer clips may be more efficient (annotation cost wise).

**Event Relation Annotation w.r.t. Middle Event:** We note there are two alternatives to our proposed annotation strategy for event relation which involves only annotating only all events only with respect to middle event. First, exhaustively annotate all event-event relations which would result in 10 annotations per video. Clearly, this is a  $2.5\times$  the annotation (in practice it is even more challenging). As a result, we decided to restrict to only one event relation. Second option is to allow choosing one of the 2-second intervals as the main event and annotating event relations with respect to it. In practice, we found the choice of main event to be subjective and inconsistent across annotations. Moreover, choosing the main event could lead to biased event relations (for instance “Caused By” relation would be more pronounced). Thus, we simplified the step by choosing Event 3 spanning from 4-6 seconds as the main event and annotated other events with respect to Event 3.

**Worker Qualification and Quality Control.** To ensure

	Acc@1		Acc@5		Recall@5	
	10 A	20 A	10 A	20 A	10 A	20 A
Majority	0.20	0.21	0.66	0.75	0.03	0.02
Human	0.62	0.71	0.96	1.00	0.64	0.59

Table 1: 10A and 20A denote 10 and 20 annotations respectively. Majority denotes choosing most frequent verbs for the validation set.

that annotators have understood the task requirements, we put up a qualification task where a worker has to successfully annotate 3 videos. These annotations are manually verified by the first author who then provides feedback on their annotations. To filter potential workers, we restrict to more than 95% approval rate and having done at least 500 tasks. In total we qualified around 120 annotators, with at least 60 workers annotating more than 30 videos every batch of  $2K$  videos.

In addition to manual qualification, we put automated checks one average number of unique verbs provided within a video, and average description lengths. We further manually inspect around 3 random samples from every annotator after every  $3K - 5K$  videos and provide constant feedback.

**Annotating Validation and Test Sets.** We ran a controlled experiments using 100 videos and annotated 25 verbs for each event. We report the human agreement in Table 1. To compute human agreement score for any event, we use one human annotation (out of 25) as a prediction and the remaining 10 or 20 annotations as ground-truths (denoted by 10A or 20A). The final score is the average over all possible prediction/ground-truth partitions. Essentially, we find that even moving from 10 to 20 annotations, the human agreement improves from 62% to 71% which suggests even at higher number of annotations, we receive verbs which are suggested by a single annotator (and hence no agreement). This rules out metrics like accuracy, precision, or F1 scores because they would penalize predictions that may be correct but are not present in a reasonably sized set of ground truth annotations. This analysis leads us to the metric Recall@5 which measures if the verbs most agreed upon by humans are indeed recalled by the model in its top-5 predictions.

Furthermore, this prompts us to collect the annotations for validation and test set in two-stages, in the first stage we collect 9 additional annotations for verb and then in the second-stage 3 annotations for argument roles and event relations given the verb (we choose the set of verbs chosen by the annotator with the highest agreement, followed by highest number of unique verbs within the video). We find this two-stage process to be of similar cost of obtaining 5 independent annotations but with the added advantage of being comparable across annotations. In total we annotation 3789 videos for validation and test sets.



	Total	Caused By	Reaction To	Enabled By	No Relation
Train Set	94016	16.94	24.05	33.76	25.25
Val Set	5304	20.99	20.29	33.82	24.88
Val Set*	4089 (77.09%)	15.3	18.95	39.05	26.66
Test Set	6392	20.19	34.88	24.44	20.4
Test Set*	4851 (75.89%)	13.39	19.04	40.9	26.5

Table 2: The distribution of Event Relations before and after filtering by taking consensus of at least two workers *i.e.* we consider only those instances where two workers agree on the event relation when given the verb.

**Reward.** We set the reward for annotating one 10-second video (for training videos) to \$0.75 after estimating the average time of completing an annotation to be around 5mins. This translates to around \$9/hour. Overall, we received generous reviews for the reward on popular turk management website. For validation and test sets, we set the reward to \$0.2 for the first stage (collecting only verbs from 9 annotators and \$0.7 for the second-stage (collecting argument and event relations from 3 annotators). As a result, the cost for annotating a single video in the validation and test set turns out to be \$3.9 ( $0.2 \times 9 + 0.7 \times 3$ ) which is around  $5.2\times$  the cost of annotating a single training video. Total cost for the process comes around \$36.7K (note: this doesn’t account for pilot experiments, qualifications, and discarded annotations due to human errors).

**Collection Timeline.** Collecting the entire training set was done over a period of about 1.2 months, and an additional 1 month for collecting the validation and test sets.

## C. Additional Dataset Statistics

In this section we report additional dataset statistics not included in Section 4.2 due to space constraints.

In Table 2 we report the distributions of Event Relations before and after filtering for validation and test sets. For filtering, we use consensus of two workers *i.e.* at least two workers agree on the argument relation which we use as the ground-truth. We largely find that the consensus on Caused By and Reaction To is low, but Enabled By and No Relations are higher.

Next, we plot the distributions for the 100 most frequent verbs, genres and chosen movies in Figure 3. For verbs and genres we find Zipf’s law in action. For verbs, we find most common verbs such as “talk”, “speak”, “walk”, “look” which are also part of frequent atomic actions despite explicitly not scoring them. This is an inherent effect due to the movie domain where dialogue is a large focus. For genres we find that “Comedy”, “Drama”, “Action”, “Romance” are the most frequent which tend to have more movements than “Mystery”, “Thriller” which have less movements on actors with often extended still-frames.

In Figure 4 we plot the top 50 most frequent words

within the argument (after removing stop-words). We find “man”, “woman” are the most frequent word in all of Arg0, Arg1, Arg2 which is not surprising since the movies are human-centric. We note the over-abundance of “man” compared to “woman” is an amplification of the biases present in the movie. Interestingly, the distribution is less skewed for Location, Direction, and Manner

## D. Implementation Details

We detail some of the implementation details for our models. All implementations are coded in PyTorch [56]. Unless otherwise mentioned we use Adam [33] optimizer with learning rate of  $1e^{-4}$ .

### D.1. Verb Prediction Models

All our implementations for verb prediction models such as I3D[8], Slow-only and SlowFast networks [16] is based on the excellent repository SlowFast [15]. We use the checkpoints from the repository for kinetics pre-trained models. All models are trained with a batch size of 8 for 10 epochs, and the model with best recall@5 is chosen for testing. For classification, we use a set of 1560 verbs composed two MLP projections (first projects to half the input dimension, the second to 1560 verbs) separated with a ReLU activation. For inference, we choose the top-5 scoring verbs. Training requires considerable GPU space, and on 8 TITAN GPUs, with batch size of 8 each epoch takes around 1 hour, with total being 10 hours.

### D.2. Argument Prediction Models

We extract the features from underlying base networks which is 2048 and 2304 for I3D and SlowFast respectively. For transformers, we use the implementation provided in Fairseq library [53]<sup>4</sup> and for GPT2 (medium) and Roberta (base) we use the implementation by HuggingFace transformer library [78]<sup>5</sup>. For tokenization and vocabulary, we utilize Byte-Pair Encoding and add special argument tokens such as [Arg0] to encode the phrases.

For both transformer encoder and decoder we use 3 layers with 8 attention heads. The decoder uses the last encoder layer outputs as encoder attention for subsequent decoding. For training, we use cross-entropy loss over the predicted sequence. For sequence generation, we use greedy-decoding with temperature 1.0 as we didn’t find improvements using beam-search or using different temperature.

For training, we used a batch size of 16 for all models other than GPT2 for which we could only use a batch size of 8 due to memory restrictions. Training time for GPT2 is around 10 hours over 8 GPUs (recall that GPT2 medium has 24 transformer layers and 16 attention heads). All other

<sup>4</sup><https://github.com/pytorch/fairseq/>

<sup>5</sup><https://github.com/huggingface/transformers>

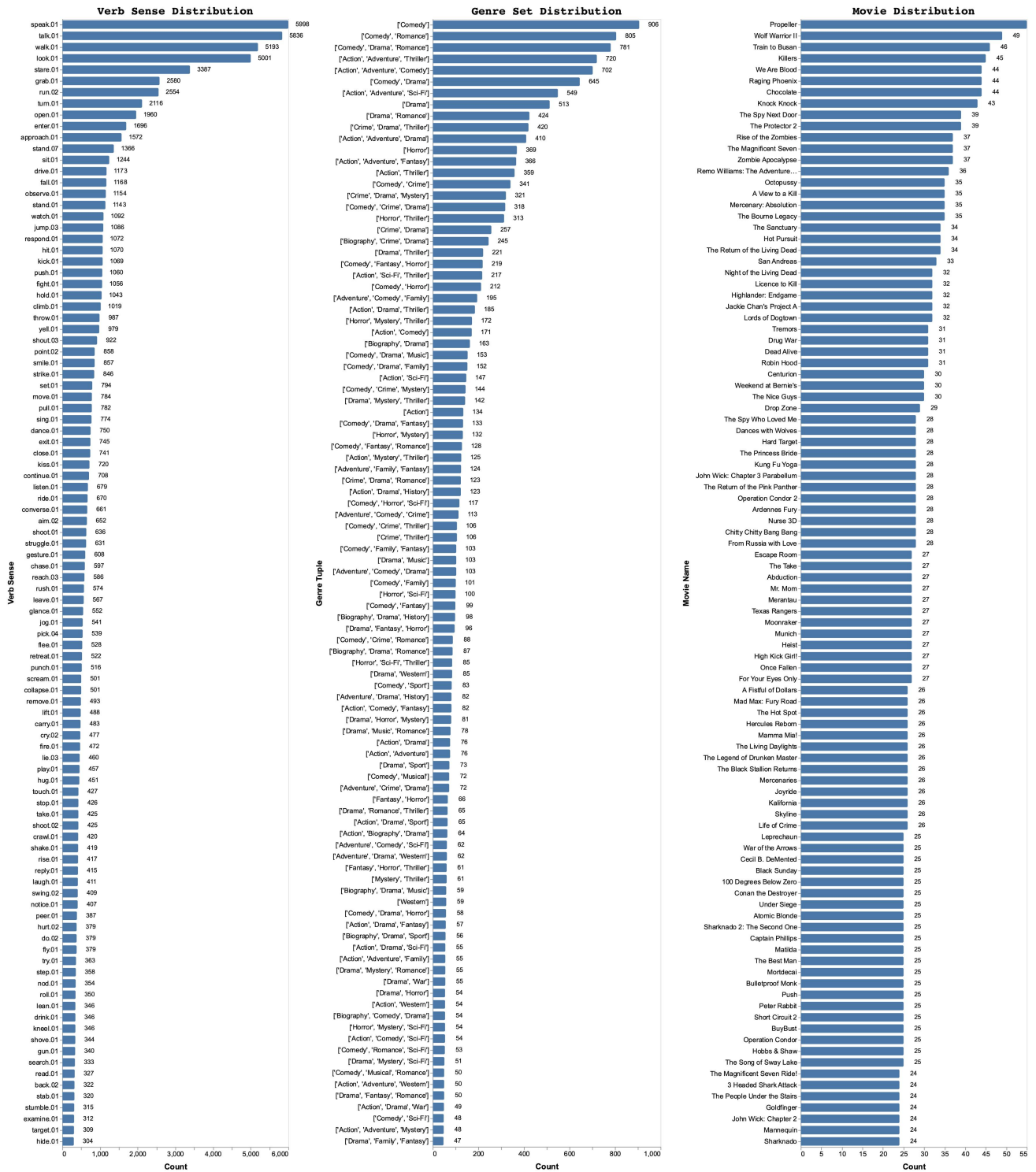


Figure 3: Distribution of 100 most frequent verbs (a), genre tuples (b), and movies (c). Note that for (a), the count represents the number of events belonging to the particular verb, whereas for (b), (c) it represents the number of video segments belonging to a particular genre or movie.

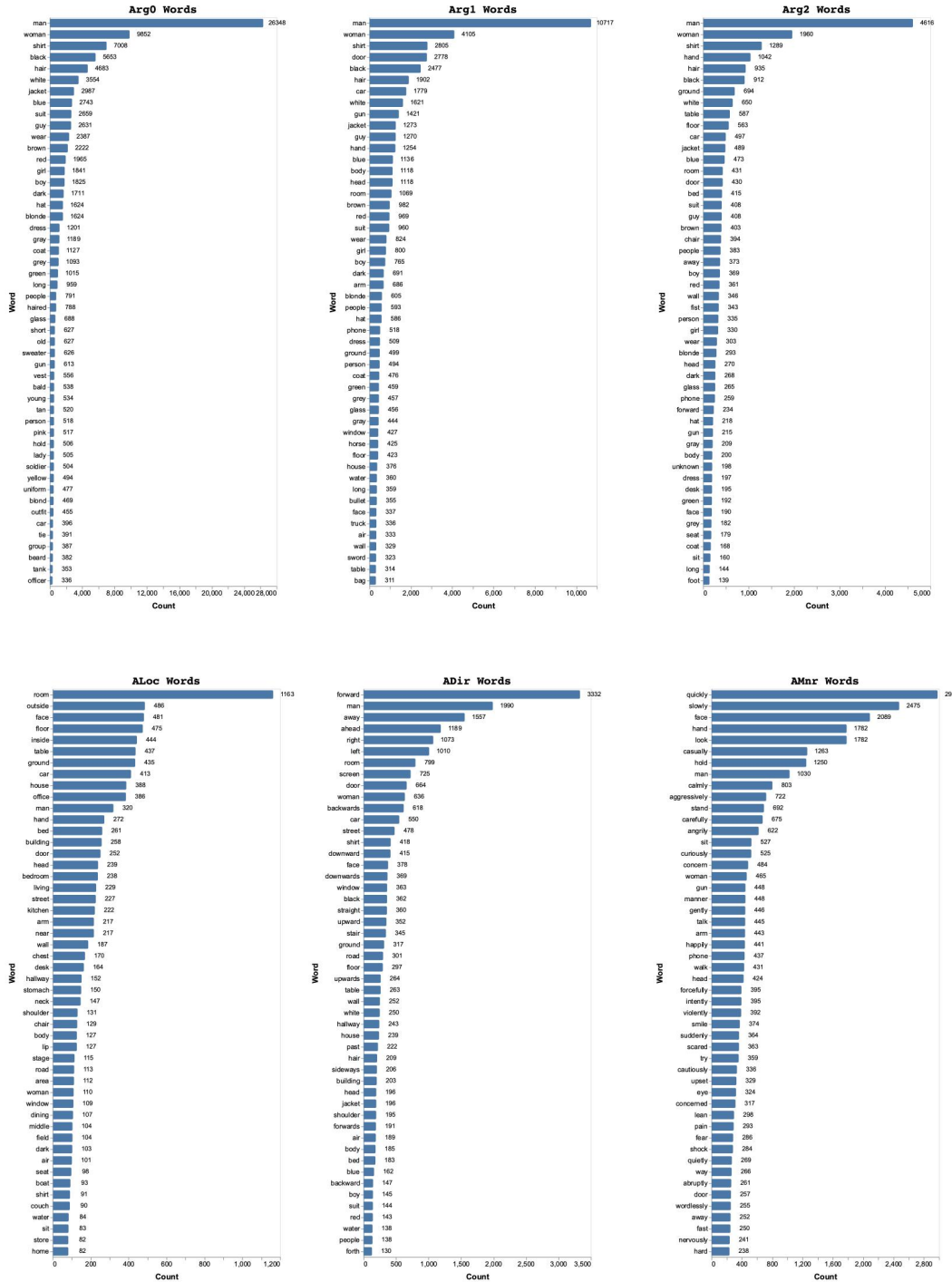


Figure 4: 50 Most frequent words (after removing stop-words) for Arg0, Arg1, Arg2, ALoc (location), ADir (direction ) and AMnr(Manner).

models take around 15 mins per epoch with batch size of 16 on a single TITAN GPU with total time around 3 hours for 10 epochs which we found sufficient for convergence.

For computing natural language generation metrics like

ROUGE, CIDEr we use the official MSCOCO Captions implementation [45] <sup>6</sup>. For co-reference metrics, we use the

<sup>6</sup><https://github.com/tylin/coco-caption>

implementation provided in coval [52]<sup>7</sup>

## E. Evaluation Metrics

In this section, we provide details on LEA as well as our proposed LEA-soft. We further report additional metrics such BLEU [55] and METEOR [5], and coreference metrics. We also report per-argument scores for the baselines.

### E.1. Co-Reference Metrics

We primarily use the metric LEA [52] which is a link-based metrics. We also note there exists other metrics such as MUC [73], BCUBE [2], CEAFE[47]. We point the reader to a seminal paper on visualizing these metrics [57] for a brief overview of MUC, BCUBE and CEAFE, and [52] for comparison of other metrics with LEA.

**LEA and LEA-soft** As noted in the paper [52], LEA computes an importance score and resolution score for each entity given as

$$\frac{\sum_{e_i \in E} \text{imp}(e_i) \times \text{res}(e_i)}{\sum_{e_i \in E} \text{imp}(e_i)} \quad (\text{E.1})$$

The final score is the F1-measure computed based on recall (entities are ground-truths) and precision (entities are predictions). As noted earlier, LEA doesn’t consider if the proposed entity by itself is correct and thus even incorrect entity predictions could lead high co-reference score as long as the co-referencing is correct. We address this using LEA-soft which additionally weights the importance of each entity during precision computation with the sum of cider scores in the numerator and len of cider scores in the denominator.

As a result, we have

$$\text{Prec}_{LEA} = \frac{\sum_{e_i \in E} \text{imp}(e_i) \times \text{res}(e_i)}{\sum_{e_i \in E} \text{imp}(e_i)} \quad (\text{E.2})$$

$$\text{Prec}_{LEA\text{-}soft} = \frac{\sum_{e_i \in E} (\sum_{e_j \in E} C(e_j)) \times \text{imp}(e_i) \times \text{res}(e_i)}{\sum_{e_i \in E} |e_i| \times \text{imp}(e_i)} \quad (\text{E.3})$$

where  $C(e_i)$  denotes the cider score for the  $i^{th}$  entity. We keep the recall computation unchanged and use the modified precision to compute the final F1-Score for LEA-soft. Since we have multiple ground-truth reference, we compute the F1-score for each ground-truth reference individually and average over the 3 ground-truths.

### E.2. Evaluation of Arguments

We examine the cider scores for different arguments over a set of 100 videos (same used for verb prediction

	cider	Arg0	Arg1	Arg2	ALoc	AScn	ADir	AMnr
GPT2	0.39	0.40	0.39	0.45	0.43	0.22	0.37	0.15
Human	0.70	0.73	0.74	0.73	0.90	0.96	0.40	0.15

Table 3: CIDER score for all collected Arguments with 5 annotations on 100 videos.

results). To compare semantic role values, which are free-form text phrases, we compute CIDER metric treating one of the chosen annotations as a hypothesis and remaining annotations as references for each argument. Table 3 compares CIDER scores for all semantic roles and scores by argument type for a GPT2 based language only baseline that generates the sequence of roles and values given the verb for an event. We find that human-agreement is high for all arguments except direction (ADir) and manner (AMnr). For both “direction” (ADir) and “manner” (AMnr), we find that both language-only baseline and human agreements are poor. On further inspection, we find that the argument “manner” describes “how” the event took place is open to subjective interpretation, and the argument “direction” has a wide range of correct values (e.g. for “walk” directions “forward”, “down the path”, and “through the trees”) may all be correct. For a reliable evaluation, we evaluate argument prediction performance only on arguments that achieved high human-agreement *i.e.* Arg0, Arg1, Arg2, ALoc, and AScn, and leave the evaluation of Direction and Manner for future work.

### E.3. All Metrics

We report BLEU@1, BLUE@2, METEOR, ROUGE, and CIDER for both val (Table 4) and test set (Table 5). For each metric we further report macro-averaged scores across verbs and arguments, and report per argument scores. Note that only CIDER is able to take advantage of the macro-averaged scores due to its inverse document frequency re-weighting. Finally, we report the co-reference metrics MUC, BCUBE, CEAFE, LEA and our proposed metric LEA-Soft.

## F. VidSitu DataSheet

The seminal work datasheets for datasets [18] outlines a list of questions to encourage transparency, accountability and mitigate unwanted biases. Here, we provide a datasheet for VidSitu closely following the guidelines in prior work. For simplicity and readability, we paste the questions verbatim.

### F.1. Motivation

- **For what purpose was the dataset created?** The main motivation to create the dataset is to bridge the re-

<sup>7</sup><https://github.com/ns-moosavi/coval>



Model Vis Feats	GPT2 X	TxDec X	Vid TxDec SlowFast	Vid TxEncDec SlowFast	Vid TxDec I3D	Vid TxEncDec I3D	Human
B@1	40.91	42.79	43.45	44.65	41.69	45.3	43.56
B@1-Vb	38.08	41.02	39.59	41.98	38.96	40.54	39.93
B@1-Arg	40.91	42.62	42.89	44.49	40.18	44.6	41.69
B@1-Arg0	44.67	46.32	48.26	48.14	49.58	49.36	49.71
B@1-Arg1	31.88	31.69	32.81	34.72	34.76	36.17	40.61
B@1-Arg2	34.13	36.3	34.93	35.86	35.17	37.36	39.87
B@1-ALoc	46.88	48.07	48.97	51.39	42.73	49.37	38.7
B@1-AScn	46.99	50.74	49.48	52.33	38.66	50.71	39.56
B@2	27.66	28.8	29.87	30.86	28.47	30.73	29.89
B@2-Vb	23.92	26.52	25.73	27.54	25	25.39	25.14
B@2-Arg	27.63	28.4	29.19	30.61	26.82	30.06	28.37
B@2-Arg0	31.06	32.07	34.09	33.78	35.33	34.03	34.74
B@2-Arg1	19.53	19.87	20.25	22.39	22.3	22.6	26.72
B@2-Arg2	22.1	23.52	22.22	23.46	21.81	24	26.76
B@2-ALoc	32.92	32.24	34.19	35.98	28.58	34.04	27.06
B@2-AScn	32.55	34.29	35.21	37.42	26.06	35.61	26.59
M	16.99	17.51	17.28	18.26	17.68	18.32	22.24
M-Vb	15.33	16.4	15.8	17.14	16.39	16.77	22.08
M-Arg	15.88	16.03	16.2	17.23	15.93	16.95	21.02
M-Arg0	21.12	21.97	20.99	21.46	22.23	22.05	25.21
M-Arg1	15.49	14.81	13.94	16.14	15.93	16.16	22.22
M-Arg2	14.99	16.27	15.21	15.65	14.76	14.85	20.75
M-ALoc	15.21	13	15.03	16.26	12.17	15.19	17.88
M-AScn	12.59	14.11	15.85	16.63	14.54	16.51	19.02
R	40.08	41.19	40.61	42.66	40.67	42.41	39.77
R-Vb	37.07	37.89	36.89	39.18	36.38	38.14	39.16
R-Arg	39.62	40.47	39.58	41.96	38.56	41.39	38.43
R-Arg0	44.77	46.7	46.78	47.36	48.65	47.71	45.84
R-Arg1	34.25	33.24	32.83	35.7	34.66	36.65	40.23
R-Arg2	33.72	36.14	34.12	35.13	34.71	35.85	36.43
R-ALoc	42.87	41.41	39.82	44.6	32.22	41.49	34.38
R-AScn	42.46	44.84	44.33	46.99	42.55	45.26	35.25
C	34.67	35.68	44.78	45.52	47.14	47.06	84.85
C-Vb	42.97	47.5	49.97	55.47	51.61	51.67	91.7
C-Arg	34.45	32.15	41.24	42.82	41.29	42.76	80.15
C-Arg0	28.33	32.1	41.64	34.6	48.99	39.42	88.24
C-Arg1	38.58	38.47	41.42	45.47	45.42	47.06	83.37
C-Arg2	36.82	40.51	42.28	41.02	40.19	44.52	74.82
C-ALoc	47.77	27.05	43.01	46.97	33.75	39.75	76.72
C-AScn	20.73	22.62	37.86	46.05	38.11	43.03	77.62
MUC	59.13	64.54	45.59	65.48	46.01	61.57	80.75
BCUBE	73.53	74.43	69.39	72.97	68.74	73.34	86.32
CEAFE	61.75	63.84	57.26	59.7	56.2	61.16	77.8
LEA	48.08	51.76	37.88	50.48	37.89	48.92	72.1
LEA Soft	28.1	28.6	28.69	31.99	30.38	33.58	70.33

Table 4: Semantic Role Prediction on Validation Set. B@1: Bleu-1, B@2: Bleu-2, M: METEOR, R: ROUGE-L, C: CIDEr, Metric-Vb: Macro Averaged over Verbs, Metric-Arg: Macro Averaged over arguments, Metric-Argi: Metric computed only for the particular argument.

Model Vis Feats	GPT2 ✗	TxDec ✗	Vid TxDec SlowFast	Vid TxEncDec SlowFast	Vid TxDec I3D	Vid TxEncDec I3D	Human
B@1	41.89	42.9	43.4	45.36	43.69	45.56	43.46
B@1-Vb	38.41	39.4	39.28	41.03	39.43	40.52	39.73
B@1-Arg	41.9	42.56	42.84	45.25	42.04	44.83	41.47
B@1-Arg0	45.65	46.06	47.56	48.92	48.96	49.75	48.2
B@1-Arg1	32.17	31.53	33.15	34.46	33.93	35.42	41.06
B@1-Arg2	35.02	37.34	34.85	36.69	36.32	38.55	39.69
B@1-ALoc	48.7	46.53	48.74	52.95	43.91	49.18	36.74
B@1-AScn	47.94	51.34	49.88	53.23	47.07	51.25	41.65
B@2	28.43	29.15	30.08	31.64	30.34	31.34	29.43
B@2-Vb	24.25	25.49	25.83	26.9	25.45	26.22	24.37
B@2-Arg	28.41	28.7	29.42	31.56	28.79	30.59	27.95
B@2-Arg0	31.69	31.92	33.56	34.33	34.84	34.76	32.99
B@2-Arg1	19.8	19.88	20.98	22.69	22.3	22.46	26.88
B@2-Arg2	22.43	24.39	22.36	24.15	23.05	24.81	26.27
B@2-ALoc	34.36	31.63	34.18	37.95	30.69	34.32	25.66
B@2-AScn	33.76	35.67	36.03	38.66	33.05	36.62	27.93
M	17.74	17.67	17.45	18.83	18.22	18.7	21.86
M-Vb	15.8	15.84	15.72	17.02	16.92	16.83	22.44
M-Arg	16.63	16.21	16.46	17.9	16.63	17.44	20.55
M-Arg0	21.82	21.83	20.72	21.96	22.2	22.23	24.61
M-Arg1	15.99	14.97	14.39	16.31	16.28	16.53	21.55
M-Arg2	15.39	16.63	15.15	16.22	15.34	15.41	20.11
M-ALoc	16.41	12.96	15.76	17.63	13.59	16.2	16.89
M-AScn	13.55	14.68	16.3	17.36	15.74	16.82	19.58
R	41.33	41.45	41.12	43.46	41.5	42.96	40.04
R-Vb	37.71	36.96	36.66	38.6	36.69	37.72	39.24
R-Arg	40.91	40.65	40.14	42.88	39.68	42.04	38.55
R-Arg0	45.89	46.6	46.75	48.22	48.69	48.3	45.5
R-Arg1	35.13	33.05	33.35	35.67	34.9	36.34	40.03
R-Arg2	34.13	36.83	33.77	35.26	35.58	36.49	37.29
R-ALoc	45.33	40.96	41.53	47.17	35.1	43.06	32.94
R-AScn	44.04	45.82	45.31	48.08	44.14	46.04	36.97
C	36.48	35.34	44.95	47.25	47.9	48.51	83.68
C-Vb	44.27	44.44	49.46	52.92	51.29	53.88	87.78
C-Arg	36.51	32.06	41.98	45.48	43.62	44.53	79.29
C-Arg0	26.17	27.83	36.84	33.51	41.89	38.64	81.62
C-Arg1	39.08	37.99	42.93	43.79	46.53	46.47	81.47
C-Arg2	35.36	41.93	39.16	39.48	41.66	43.84	73.21
C-ALoc	55.05	25.83	48.3	58.38	43.83	45.15	77.38
C-AScn	26.9	26.71	42.65	52.22	44.18	48.57	82.77
MUC	60.51	65.42	47.51	65.91	47.63	62.62	80.8
BCUBE	74.21	74.76	69.84	72.95	69.2	73.6	86.26
CEAFE	62.19	63.85	57.33	59.57	56.65	61.41	77.38
LEA	49.38	52.46	38.91	50.88	38.77	49.61	71.77
LEA Soft	30.24	29.18	30.21	33.5	31.73	35.46	70.6

Table 5: Semantic Role Prediction on Test Set. B@1: Bleu-1, B@2: Bleu-2, M: METEOR, R: ROUGE-L, C: CIDEr, Metric-Vb: Macro Averaged over Verbs, Metric-Arg: Macro Averaged over arguments, Metric-Argi: Metric computed only for the particular argument.

search gap between learning atomic actions and generating holistic captions. In particular, the dataset opens path for the task of Visual Semantic Role Labeling in Videos which in addition to action-recognition, emphasizes how various objects interact within an action, how various objects interact over time-period across multiple actions, co-referencing of these objects over time, and how various actions affect each other.

- **Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?** The dataset is created by the authors who belong to PRIOR Team at AI2. The first-author (Arka Sadhu) was a summer intern in the PRIOR Team.
- **Who funded the creation of the dataset?** PRIOR Team at AI2 funded the creation of the dataset.

## F.2. Composition

- **What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?** Each instance consists of a 10-second video obtained from a movie-clip available on YouTube. These are usually human-centric and hence primarily contain videos of people interacting in diverse and complex situations.
- **How many instances are there in total (of each type, if appropriate)?** In total there are 27.4K instances distributed across training (23.62K), validation (1.80K) and testing (1.98K)
- **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?** This question doesn't pertain to our dataset.
- **What data does each instance consist of?** Each instance is a 10-second video (mp4 video) available from YouTube.
- **Is there a label or target associated with each instance?** Each instance (10 second video) is annotated at 2-second intervals with a verb describing the event, corresponding argument roles for the verb co-referenced across the video, and event relations across the various verbs with respect to the middle event (Event 3 spanning from 4-6 seconds).
- **Is any information missing from individual instances?** No, every instance has the same annotations.
- **Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)?** We provide information about which instances are derived from the same 2 – 3 minutes YouTube video as well as the underlying movie (this information is obtained from Condensed-Movies [3]

dataset). However, this information is not used for any of the task in the dataset except for splitting the videos in train, validation and test sets.

- **Are there recommended data splits (e.g., training, development/validation, testing)?** Yes, we provide training, validation and test sets by splitting the overall set in 80 : 10 : 10 ratio randomly based on the movie names. We also ensure (qualitatively) that the normalized distributions of verbs, and genres are same across the splits.
- **Are there any errors, sources of noise, or redundancies in the dataset?** The main sources of errors would be the annotations themselves, however, we have made extended efforts from automatic to manual checks to remove such errors and provided constant feedback. Some redundancy may occur due to oversampling of dialogues in movies which are described with the verb "talk". Some redundancy may also occur due to use of closely related verbs such as "run" and "jog".
- **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)?** Yes, the dataset provides links to YouTube videos. Since the videos are provided by a licensed channel, we expect the videos to have high online longevity.
- **Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor patient confidentiality, data that includes the content of individuals' non-public communications)?** No, our dataset is derived from movies publicly available on youtube.
- **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?** Some of the videos obtained from action, crime or horror movies may be sensitive to some viewers when viewed directly. Some videos may also contain violence and gore, and we suggest user discretion in viewing the videos.

## F.3. Collection Process

- **How was the data associated with each instance acquired?** The data was directly observable in the form of embedded youtube videos.
- **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)?** We used Amazon Mechanical Turk to collect the data with a custom annotation interface. We validated them by small scale user study and taking feedbacks during worker qualification.

- **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?** We sampled videos which had more verbs within their duration.
- **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?** Crowd-Workers were involved in the process. They were paid \$0.75 for training videos and \$0.2 for verb annotation and \$0.7 for argument and event relation for videos in validation and test splits. On average it is around \$9 – \$12 per hour above the minimum wage. On popular websites, our pay was noted to be generous.
- **Over what timeframe was the data collected?** The data was collected over 2.2 months with initial 1.2 months for training set and rest for validation and testing.
- **Were any ethical review processes conducted (e.g., by an institutional review board)?** No, there was no ethical review process.

#### F.4. Preprocessing/cleaning/labeling

- **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?** Only, exact string match was performed to obtain co-referenced entities. We used spacy [27] to compute dataset statistics such as noun-diversity but it is not used over the collected data for down-stream tasks.
- **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)?** In our case, raw data is same as cleaned data.

#### F.5. Uses

- **Has the dataset been used for any tasks already?** We have used the data to show its usefulness for our proposed task Visual Semantic Role Labeling in Videos
- **Is there a repository that links to any or all papers or systems that use the dataset?** Updated information about the dataset can be found on [vidsitu.org](https://vidsitu.org).
- **What (other) tasks could the dataset be used for?** We believe the dataset could be re-purposed for many down-stream video understanding tasks such as video retrieval, video question answering, action forecasting, long-term reasoning.

- **Are there tasks for which the dataset should not be used?** The data is obtained from movies and exhibits certain stereotypes which donot hold true in real world. It also contains highly unlikely action sequences (such as a “man flying”), and thus it shouldn’t be used for real-world cases and strictly used as a video understanding benchmark.

#### F.6. Distribution

- **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?** The dataset is publicly available at [vidsitu.org](https://vidsitu.org).
- **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** The dataset is available through our website and github. The dataset is stored on Amazon S3 buckets.



## References

- [1] Nayyer Aafaq, Ajmal Mian, Wei Liu, Syed Zulqarnain Gilani, and Mubarak Shah. Video description: A survey of methods, datasets, and evaluation metrics. *ACM Computing Surveys (CSUR)*, 52(6):1–37, 2019. 3
- [2] A. Bagga and B. Baldwin. Entity-based cross-document coreferencing using the vector space model. In *COLING-ACL*, 1998. 7, 16
- [3] M. Bain, Arsha Nagrani, A. Brown, and Andrew Zisserman. Condensed movies: Story based retrieval with contextual embeddings. In *ACCV*, 2020. 3, 4, 9, 19
- [4] C. Baker, C. Fillmore, and J. Lowe. The berkeley framenet project. In *COLING-ACL*, 1998. 9
- [5] S. Banerjee and A. Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *IEEevaluation@ACL*, 2005. 16
- [6] Claire Bonial, Olga Babko-Malaya, Jinho D Choi, Jena Hwang, and Martha Palmer. Propbank annotation guidelines. *Center for Computational Language and Education Research, CU-Boulder*, 2010. 9
- [7] Susan Windisch Brown, Julia Bonn, James Gung, Annie Zainen, James Pustejovsky, and Martha Palmer. VerbNet representations: Subevent semantics for transfer verbs. In *Proceedings of the First International Workshop on Designing Meaning Representations*, pages 154–163, Florence, Italy, Aug. 2019. Association for Computational Linguistics. 9
- [8] J. Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4724–4733, 2017. 2, 6, 13
- [9] Yu-Wei Chao, Yunfan Liu, Xieyang Liu, Huayi Zeng, and Jia Deng. Learning to detect human-object interactions. *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 381–389, 2018. 3
- [10] Yu-Wei Chao, Z. Wang, Yugeng He, J. Wang, and Jia Deng. Hico: A benchmark for recognizing human-object interactions in images. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1017–1025, 2015. 3
- [11] Zhenfang Chen, L. Ma, Wenhan Luo, and K. Wong. Weakly-supervised spatio-temporally grounding natural sentence in video. In *ACL*, 2019. 3
- [12] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The epic-kitchens dataset. In *European Conference on Computer Vision (ECCV)*, 2018. 3
- [13] P. Das, C. Xu, R. F. Doell, and Corso J. J. A thousand frames in just a few words: Lingual description of videos through latent topics and sparse object stitching. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2013. 3
- [14] Ali Diba, Mohsen Fayyaz, Vivek Sharma, Manohar Paluri, Jürgen Gall, Rainer Stiefelhofen, and Luc Van Gool. Large scale holistic video understanding. In *European Conference on Computer Vision*, pages 593–610. Springer, 2020. 3
- [15] Haoqi Fan, Yanghao Li, Bo Xiong, Wan-Yen Lo, and Christoph Feichtenhofer. Pyslowfast. <https://github.com/facebookresearch/slowfast>, 2020. 13
- [16] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6201–6210, 2019. 2, 6, 10, 13
- [17] J. Gao, C. Sun, Zhenheng Yang, and R. Nevatia. Tall: Temporal activity localization via language query. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5277–5285, 2017. 3
- [18] Timnit Gebru, J. Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, H. Wallach, Hal Daumé, and K. Crawford. Datasheets for datasets. *ArXiv*, abs/1803.09010, 2018. 9, 16
- [19] Rohit Girdhar, J. Carreira, C. Doersch, and Andrew Zisserman. Video action transformer network. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 244–253, 2019. 2
- [20] Raghav Goyal, S. Kahou, Vincent Michalski, Joanna Materzynska, S. Westphal, Heuna Kim, V. Haenel, Ingo Fründ, Peter Yianilos, Moritz Mueller-Freitag, F. Hoppe, Christian Thureau, I. Bax, and R. Memisevic. The “something something” video database for learning and evaluating visual common sense. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5843–5851, 2017. 3
- [21] C. Gu, C. Sun, Sudheendra Vijayanarasimhan, C. Pantofaru, D. Ross, G. Toderici, Y. Li, S. Ricco, R. Sukthankar, C. Schmid, and J. Malik. Ava: A video dataset of spatio-temporally localized atomic visual actions. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6047–6056, 2018. 2, 3, 9, 10
- [22] Saurabh Gupta and Jitendra Malik. Visual semantic role labeling. *ArXiv*, abs/1505.04474, 2015. 2, 3
- [23] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask r-cnn. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42:386–397, 2020. 10
- [24] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Nibbles. Activitynet: A large-scale video benchmark for human activity understanding. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 961–970, 2015. 2, 3, 9
- [25] Lisa Anne Hendricks, O. Wang, E. Shechtman, Josef Sivic, Trevor Darrell, and Bryan C. Russell. Localizing moments in video with natural language. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5804–5813, 2017. 3
- [26] Y. Hong, Tongtao Zhang, Timothy J. O’Gorman, Sharone Horowitz-Hendler, Huai zhong Ji, and Martha Palmer. Building a cross-document event-event relation corpus. In *LAW@ACL*, 2016. 4, 10
- [27] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017. 20
- [28] Qingqiu Huang, Yu Xiong, Anyi Rao, Jiase Wang, and Dahua Lin. Movienet: A holistic dataset for movie under-

- standing. In *The European Conference on Computer Vision (ECCV)*, 2020. 3, 9
- [29] Haroon Idrees, Amir R. Zamir, Yu-Gang Jiang, Alex Gorbun, Ivan Laptev, Rahul Sukthankar, and Mubarak Shah. The thumos challenge on action recognition for videos “in the wild”. *Computer Vision and Image Understanding*, 155:1–23, 2017. 3
- [30] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black. Towards understanding action recognition. In *International Conf. on Computer Vision (ICCV)*, pages 3192–3199, Dec. 2013. 3
- [31] W. Kay, J. Carreira, K. Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, F. Viola, T. Green, T. Back, A. Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *ArXiv*, abs/1705.06950, 2017. 2, 3, 6, 9
- [32] Mert Kilickaya, Aykut Erdem, Nazli Ikizler-Cinbis, and Erkut Erdem. Re-evaluating automatic metrics for image captioning. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 199–209, Valencia, Spain, Apr. 2017. Association for Computational Linguistics. 2
- [33] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015. 13
- [34] Y. Kong and Yun Fu. Human action recognition and prediction: A survey. *ArXiv*, abs/1806.11230, 2018. 3
- [35] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *International Conference on Computer Vision (ICCV)*, 2017. 2, 5, 12
- [36] Hilde Kuehne, Hueihan Jhuang, E. Garrote, T. Poggio, and Thomas Serre. Hmdb: A large video database for human motion recognition. *2011 International Conference on Computer Vision*, pages 2556–2563, 2011. 3
- [37] Jie Lei, Licheng Yu, Mohit Bansal, and T. Berg. Tvqa: Localized, compositional video question answering. In *EMNLP*, 2018. 3
- [38] Jie Lei, Licheng Yu, T. Berg, and Mohit Bansal. Tvr: A large-scale dataset for video-subtitle moment retrieval. In *ECCV*, 2020. 3
- [39] Ang Li, Meghana Thotakuri, D. Ross, J. Carreira, Alexander Votrikov, and Andrew Zisserman. The ava-kinetics localized human actions video dataset. *ArXiv*, abs/2005.00214, 2020. 3
- [40] Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. Hero: Hierarchical encoder for video+ language omni-representation pre-training. In *EMNLP*, 2020. 3
- [41] Manling Li, Alireza Zareian, Q. Zeng, Spencer Whitehead, Di Lu, Huai zhong Ji, and Shih-Fu Chang. Cross-media structured common space for multimedia event extraction. In *ACL*, 2020. 3
- [42] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. 7
- [43] T. Lin, X. Liu, Xin Li, E. Ding, and Shilei Wen. Bmn: Boundary-matching network for temporal action proposal generation. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3888–3897, 2019. 2
- [44] Tianwei Lin, Xu Zhao, Haisheng Su, Chongjing Wang, and Ming Yang. Bsn: Boundary sensitive network for temporal action proposal generation. In *European Conference on Computer Vision*, 2018. 2
- [45] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 10, 15
- [46] Y. Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, M. Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pre-training approach. *ArXiv*, abs/1907.11692, 2019. 6, 7
- [47] Xiaoqiang Luo. On coreference resolution performance metrics. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 25–32, Vancouver, British Columbia, Canada, Oct. 2005. Association for Computational Linguistics. 7, 16
- [48] Louis Mahon, Eleonora Giunchiglia, B. Li, and Thomas Lukasiewicz. Knowledge graph extraction from videos. *2020 19th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 25–32, 2020. 2
- [49] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-End Learning of Visual Representations from Uncurated Instructional Videos. In *CVPR*, 2020. 3
- [50] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In *ICCV*, 2019. 2, 3, 9
- [51] Mathew Monfort, B. Zhou, Sarah Adel Bargal, Alex Andonian, Tom Yan, K. Ramakrishnan, L. Brown, Quanfu Fan, Dan Gutfreund, Carl Vondrick, and A. Oliva. Moments in time dataset: One million videos for event understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42:502–508, 2020. 3, 9
- [52] N. Moosavi and Michael Strube. Which coreference evaluation metric do you trust? a proposal for a link-based entity aware metric. In *ACL*, 2016. 7, 16
- [53] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019. 13
- [54] Martha Palmer, Paul Kingsbury, and Daniel Gildea. The proposition bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31:71–106, 2005. 2, 3, 4, 9, 10
- [55] Kishore Papineni, S. Roukos, T. Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, 2002. 16
- [56] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming

- Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems* 32, pages 8024–8035. Curran Associates, Inc., 2019. [13](#)
- [57] Sameer Pradhan, Xiaoqiang Luo, Marta Recasens, Eduard Hovy, Vincent Ng, and Michael Strube. Scoring coreference partitions of predicted mentions: A reference implementation. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 30–35, Baltimore, Maryland, June 2014. Association for Computational Linguistics. [16](#)
- [58] Sarah Pratt, Mark Yatskar, Luca Weihs, Ali Farhadi, and Aniruddha Kembhavi. Grounded situation recognition. In *ECCV*, 2020. [2, 3](#)
- [59] Michaela Regneri, Marcus Rohrbach, Dominikus Wetzels, Stefan Thater, B. Schiele, and Manfred Pinkal. Grounding action descriptions in videos. *Transactions of the Association for Computational Linguistics*, 1:25–36, 2013. [3](#)
- [60] Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and B. Schiele. A dataset for movie description. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3202–3212, 2015. [3, 4, 5, 9, 10](#)
- [61] Arka Sadhu, K. Chen, and R. Nevatia. Video object grounding using semantic roles in language description. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10414–10424, 2020. [2, 3, 9](#)
- [62] Peng Shi and Jimmy Lin. Simple bert models for relation extraction and semantic role labeling. *ArXiv*, abs/1904.05255, 2019. [4, 5, 10](#)
- [63] Gunnar A. Sigurdsson, G. Varol, X. Wang, Ali Farhadi, I. Laptev, and A. Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *ECCV*, 2016. [3](#)
- [64] Carina Silberer and Manfred Pinkal. Grounding semantic roles in images. In *EMNLP*, 2018. [3](#)
- [65] K. Soomro, A. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *ArXiv*, abs/1212.0402, 2012. [3](#)
- [66] Emma Strubell, Patrick Verga, Daniel Andor, David Weiss, and Andrew McCallum. Linguistically-informed self-attention for semantic role labeling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5027–5038, Brussels, Belgium, Oct.-Nov. 2018. Association for Computational Linguistics. [9](#)
- [67] Chen Sun, Abhinav Shrivastava, Carl Vondrick, Kevin Murphy, Rahul Sukthankar, and Cordelia Schmid. Actor-centric relation network. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer Vision – ECCV 2018*, pages 335–351, Cham, 2018. Springer International Publishing. [2](#)
- [68] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. [6](#)
- [69] Yansong Tang, Dajun Ding, Yongming Rao, Y. Zheng, Danyang Zhang, L. Zhao, Jiwen Lu, and J. Zhou. Coin: A large-scale dataset for comprehensive instructional video analysis. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1207–1216, 2019. [3, 9](#)
- [70] Makarand Tapaswi, Y. Zhu, R. Stiefelhausen, A. Torralba, R. Urtasun, and S. Fidler. Movieqa: Understanding stories in movies through question-answering. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4631–4640, 2016. [3](#)
- [71] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. [6, 7](#)
- [72] Paul Vicol, Makarand Tapaswi, Lluís Castrejon, and Sanja Fidler. Moviegraphs: Towards understanding human-centric situations from videos. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [3, 9](#)
- [73] Marc B. Vilain, J. Burger, J. Aberdeen, D. Connolly, and L. Hirschman. A model-theoretic coreference scoring scheme. In *MUC*, 1995. [7, 16](#)
- [74] Oriol Vinyals, A. Toshev, Samy Bengio, and D. Erhan. Show and tell: Lessons learned from the 2015 mscoco image captioning challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39:652–663, 2017. [2](#)
- [75] L. Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, D. Lin, X. Tang, and L. Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. [2](#)
- [76] X. Wang, Ross B. Girshick, A. Gupta, and Kaiming He. Non-local neural networks. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7794–7803, 2018. [6](#)
- [77] Xin Eric Wang, Jiawei Wu, Junkun Chen, Lei Li, Y. Wang, and William Yang Wang. Vatec: A large-scale, high-quality multilingual dataset for video-and-language research. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4580–4590, 2019. [2, 3, 5](#)
- [78] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, Oct. 2020. Association for Computational Linguistics. [13](#)
- [79] Chao-Yuan Wu, Christoph Feichtenhofer, Haoqi Fan, Kaiming He, Philipp Krähenbühl, and Ross B. Girshick. Long-

- term feature banks for detailed video understanding. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 284–293, 2019. 2
- [80] Yu Xiong, Qingqiu Huang, Lingfeng Guo, Hang Zhou, Bolei Zhou, and Dahua Lin. A graph-based framework to bridge movies and synopses. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 9
- [81] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *ACM Multimedia*, 2017. 2, 3
- [82] J. Xu, T. Mei, Ting Yao, and Y. Rui. Msr-vtt: A large video description dataset for bridging video and language. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5288–5296, 2016. 2, 3, 5
- [83] Mark Yatskar, Luke Zettlemoyer, and Ali Farhadi. Situation recognition: Visual semantic role labeling for image understanding. In *Conference on Computer Vision and Pattern Recognition*, 2016. 2, 3, 9
- [84] Sarah Young, Dhiraj Gandhi, Shubham Tulsiani, Abhinav Gupta, P. Abbeel, and Lerrel Pinto. Visual imitation made easy. *ArXiv*, abs/2008.04899, 2020. 2
- [85] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *AAAI*, pages 9127–9134, 2019. 2, 3
- [86] Kuo-Hao Zeng, Tseng-Hung Chen, Ching-Yao Chuang, Yuan-Hong Liao, Juan Carlos Niebles, and Min Sun. Leveraging video descriptions to learn video question answering. In *AAAI*, 2017. 3
- [87] H. Zhang, Yi-Xiang Zhang, B. Zhong, Qing Lei, Lijie Yang, Ji-Xiang Du, and Duan-Sheng Chen. A comprehensive survey of vision-based human action recognition methods. *Sensors (Basel, Switzerland)*, 19, 2019. 3, 9
- [88] Zixing Zhang, Zhou Zhao, Yang Zhao, Q. Wang, H. Liu, and L. Gao. Where does it exist: Spatio-temporal video grounding for multi-form sentences. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10665–10674, 2020. 3
- [89] Hang Zhao, Antonio Torralba, Lorenzo Torresani, and Zhicheng Yan. Hacs: Human action clips and segments dataset for recognition and temporal localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8668–8678, 2019. 3
- [90] Luowei Zhou, Nathan Louis, and Jason J Corso. Weakly-supervised video object grounding from text by loss weighting and object interaction. In *British Machine Vision Conference*, 2018. 3, 9
- [91] Luowei Zhou, Yingbo Zhou, Jason J. Corso, R. Socher, and Caiming Xiong. End-to-end dense video captioning with masked transformer. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8739–8748, 2018. 2
- [92] Linchao Zhu and Y. Yang. Actbert: Learning global-local video-text representations. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8743–8752, 2020. 3