ANetQA: A Large-scale Benchmark for Fine-grained Compositional Reasoning over Untrimmed Videos

Zhou Yu¹ Lixiang Zheng¹ Zhou Zhao² Fei Wu² Jianping Fan^{1,3} Kui Ren⁴ Jun Yu^{1*}

¹ School of Computer Science, Hangzhou Dianzi University, China.

²Colledge of Computer Science and Technology, Zhejiang University, China

³AI Lab at Lenovo Research, China

⁴School of Cyber Science and Technology, Zhejiang University, China

{yuz, lxzheng, yujun}@hdu.edu.cn, {zhaozhou, wufei, kuiren}@zju.edu.cn, jfan1@lenovo.com

Abstract

Building benchmarks to systemically analyze different capabilities of video question answering (VideoQA) models is challenging yet crucial. Existing benchmarks often use non-compositional simple questions and suffer from language biases, making it difficult to diagnose model weaknesses incisively. A recent benchmark AGQA [8] poses a promising paradigm to generate QA pairs automatically from pre-annotated scene graphs, enabling it to measure diverse reasoning abilities with granular control. However, its questions have limitations in reasoning about the finegrained semantics in videos as such information is absent in its scene graphs. To this end, we present ANetQA, a large-scale benchmark that supports fine-grained compositional reasoning over the challenging untrimmed videos from ActivityNet [4]. Similar to AGQA, the QA pairs in ANetQA are automatically generated from annotated video scene graphs. The fine-grained properties of ANetQA are reflected in the following: (i) untrimmed videos with fine-grained semantics; (ii) spatio-temporal scene graphs with fine-grained taxonomies; and (iii) diverse questions generated from fine-grained templates. ANetQA attains 1.4 billion unbalanced and 13.4 million balanced QA pairs, which is an order of magnitude larger than AGQA with a similar number of videos. Comprehensive experiments are performed for state-of-the-art methods. The best model achieves 44.5% accuracy while human performance tops out at 84.5%, leaving sufficient room for improvement.

1. Introduction

Recent advances in deep learning have enabled machines to tackle complicated video-language tasks that involve



Figure 1. Comparisons of ANetQA and AGQA [8]. The QA pairs in both benchmarks are automatically generated from spatiotemporal scene graphs by using handcrafted question templates. Benefiting from the untrimmed long videos and fine-grained scene graphs, our questions require more fine-grained reasoning abilities than those in AGQA when similar templates are applied. Moreover, the newly introduced attribute annotations allow us to design many fine-grained question templates that are not supported in AGQA (*e.g., "what color"* and *"what is the occupation"*).

both video and language clues, *e.g.*, video-text retrieval, video captioning, video temporal grounding, and video question answering. Among these tasks, video question answering (VideoQA) is one of the most challenging tasks as it verifies multiple skills simultaneously. Taking the question "What is the black object that the person is wearing before various fish are seen swimming through the reef?" in Figure 1 as an example, it requires a synergistic understanding of both the video and question, together with spatio-temporal reasoning to predict an accurate answer.

^{*}Jun Yu is the corresponding author

To comprehensively evaluate the capabilities of existing VideoQA models, several prominent benchmarks have been established [11, 21, 29, 33, 38, 42, 43]. Despite their usefulness, they also have distinct shortcomings. Some benchmarks use simulated environments to synthesize video contents [29, 42], which provides controllable diagnostics over different reasoning skills. However, the synthetic videos lack visual diversity and the learned models on the benchmarks cannot generalize to real-world scenarios directly. Some real-world benchmarks generate QA pairs from off-the-shelf video captions [38, 48] or human annotations [11, 21, 33, 43], which suffer from simple question expressions and biased answer distributions. These weaknesses may be exploited by models to make educated guesses to obtain the correct answers without seeing video contents [24, 40].

One recent VideoQA benchmark AGQA poses a promising paradigm to address the above limitations [8]. AGQA is built upon the real-world videos from Charades [32]. In contrast to previous benchmarks, AGQA adopts a twostage paradigm instead. For each video, a spatio-temporal scene graph over representative frames is first annotated by humans, which consists of spatially-grounded objectrelationship triplets and temporally-grounded actions. After that, different types of questions are generated on top of the scene graph using corresponding question templates, enabling it to measure various reasoning abilities with granular control. Despite the comprehensiveness of AGQA, we argue that its foundation-the spatio-temporal scene graph—has limitations in representing the *fine-grained* semantics of videos. Specifically, their scene graphs encode objects and relationships from limited taxonomies, which are not fine-grained enough for generating questions that require reasoning about the detailed video semantics.

To this end, we introduce ANetQA¹, a new benchmark that supports fine-grained compositional reasoning over complex web videos from ActivityNet [4]. Similar to the strategy of AGQA, the QA pairs in ANetQA are automatically generated from pre-annotated scene graphs. As shown in Figure 1, we claim that ANetQA is more fine-grained than AGQA in terms of the following:

- (i) The benchmark is built upon untrimmed long videos with fine-grained semantics. Each video may involves multiple indoor or outdoor scenarios, containing complicated interactions between persons and objects.
- (ii) The spatio-temporal scene graph consists of finegrained objects (*e.g.*, "manta ray", "diving gear"), relationships (*e.g.*, "jumping into", "chasing"), attributes (*e.g.*, "swimming", "black and white"), and actions in natural language (*e.g.*, "a manta ray swims in the ocean over a reef").

(iii) Benefiting from the fine-grained scene graphs, we are able to design diverse question templates that requires fine-grained compositional reasoning (*e.g.*, *"what color ..."* and *"what is the occupation ..."*).

Benefiting from the above fine-grained characteristics, ANetQA obtains 1.4B unbalanced and 13.4M balanced QA pairs. To the best of our knowledge, ANetQA is the largest VideoQA benchmark in terms of the number of questions. Compared with the previous largest benchmark AGQA, ANetQA is an order of magnitude larger than it with a similar number of videos. We conduct comprehensive experiments and intensive analyses on ANetQA for the state-of-the-art VideoQA models, including HCRN [19], ClipBERT [20], and All-in-One [35]. The best model delivers 44.5% accuracy while human performance tops out at 84.5%, showing sufficient room for future improvement. The benchmark is available at here².

2. Related Work

We briefly review the field of VideoQA in terms of methods and benchmarks. Since ANetQA is built upon ActivityNet [4], we introduce ActivityNet and its derived benchmarks in particular.

VideoQA approaches. The research of visual question answering lies mainly in the image domain. A number of image question answering (ImageQA) methods have been developed to push state-of-the-art performance on public benchmarks successively [6, 13, 44, 45]. As a natural extension of the ImageQA task, VideoQA is more challenging as it requires effective temporal representation modeling and spatio-temporal reasoning. Existing studies explore end-to-end neural networks in conjunction with hierarchical representations [38, 49], memory networks [7, 30, 33], and graph networks [9, 25, 37].

Motivated by the encouraging success of Transformers [34] in various NLP [15, 31], CV [3, 26], and multimodal tasks [1,2,27], Transformer-based approaches have become the mainstream of recent VideoQA research. Early approaches only exploit the Transformer architecture and train models from scratch [14, 23]. More recently, pretrained Transformer models on large-scale datasets have shown effectiveness when finetuned on VideoQA tasks. Some approaches incorporate the pretrained language Transformers [16,41] or multimodal Transformers on image-text pairs [20] to improve VideoQA performance. Some other studies perform video-language pretraining directly on massive video-text pairs, which learn better multimodal representations and achieve state-of-the-art performance on various VideoQA benchmarks [5, 35, 39, 47].

VideoQA benchmarks. The rapid progress in VideoQA is inextricably related to the established benchmarks. Existing

¹Note that there is a VideoQA benchmark ActivityNet-QA [43] whose QA pairs are fully annotated by humans. To avoid confusion, we name our benchmark ANetQA.

²https://milvlg.github.io/anetqa

	video			question		grounding taxonomy			
	type	#videos	avg. len.	#QA pairs	#templates	#objects	#relations	#attributes	#actions
CLEVRER [42]	synth.	20K	5s	305K	5	1	2	13	3
TVQA+ [22]	real	4.2K	7.2s	29.4K	-	2,527	-	-	open
HowtoVQA69M [39]	real	69M	12.1s	69M	-	-	-	-	open
AGQA [8]	real	9.6K	30s	192M/3.9M	28	36	44	-	157
ANetQA	real	11.5K	180s	1.4B/13.4M	119	2,072	86	618	open

Table 1. **Comparisons of ANetQA and other representative large-scale VideoQA benchmarks**. Benefiting from the fine-grained video and grounding annotations, ANetQA attains massive fine-grained questions and is an order of magnitude larger than the current largest benchmarks [8,39] in terms of the number of QA pairs. "open" indicates the grounded actions are depicted in natural language.

VideoQA benchmarks can be categorized into two groups based on whether their videos are synthesized by simulation [29, 42] or collected from the real world [8, 17, 21, 28, 33, 36, 38, 39, 43, 46, 48]. The synthesized benchmarks can easily obtain massive QA pairs without human annotations. Their synthetic nature also enables granular control over reasoning abilities and language biases. However, the synthesized videos are often short and lack visual diversity, making it difficult to generalize the learned models to realworld scenarios. Establishing VideoQA benchmarks on real-world videos requires human annotations inevitably. Early benchmarks rely on the associated video captions to generate QA pairs automatically [28, 38, 48, 51]. Although these captions are annotated by humans, they are often too general to cover all the fine-grained semantics in videos. This makes these benchmarks be dominated by simple questions that lack detailed information. To obtain fine-grained and diverse questions, some recent benchmarks have been established by asking annotators to design questions of specific reasoning abilities, e.g., object localization [22], relationship recognition [43], and causality analysis [36]. Nevertheless, prohibitive annotation costs restrict the sizes of these benchmarks and free-form question expressions lead to severe language biases. One recent benchmark AGQA introduces a new paradigm to automatically generate QA pairs upon video scene graphs [8]. Through the composition of scene graph elements, AGQA is orders of magnitude larger than its counterparts. Similar to AGQA, our ANetQA is also built upon spatiotemporal scene graphs. In contrast to AGQA, ANetQA shows its fine-grained characteristics in terms of the videos, annotated scene graphs, and generated questions. Detailed comparisons of ANetQA and other representative largescale VideoOA benchmarks are shown in Table 1.

ActivityNet and its derivatives. ActivityNet (*abbr*: ANet) is one of the most important video recognition benchmarks [4]. It consists of 20K untrimmed videos from 200 activity classes, including both indoor and outdoor scenarios. The benchmark is challenging as its videos contain rich semantics. Therefore, some derived benchmarks are built upon ANet to provide fine-grained annotations [18, 50]. ANet-Captions [18] annotates each video with multiple

temporally-grounded captions. ANet-Entities [50] provides spatially-grounded bounding boxes for the noun phrases mentioned in the captions. We establish our ANetQA based on the annotations of these two benchmarks .

3. The ANetQA Benchmark

ANetQA is a large-scale VideoQA benchmark to measure a variety of spatio-temporal reasoning abilities at a fine-grained level. In this section, we first provide an overview of the construction process of our benchmark and then introduce the key stages in detail.

3.1. Overview

The videos in ANetQA are derived from ActivityNet [4]. As mentioned above, we leverage the auxiliary annotations on ActivityNet [18,50] to reduce the annotation costs during the construction of our benchmark. These result in 11,525 videos in total, which are comprised of 9,155 and 2,370 videos in the train and val subsets of ActivityNet, respectively. We keep the train subset unchanged and further divide the val subset evenly into a new val subset of 1,185 videos.

Next, we annotate each video with a spatio-temporal scene graph via crowdsourcing. Each video has been annotated with temporal-grounded captions [18] and spatiallygrounded objects from a few representative frames [50], For each frame, we first clean the mislabeled objects and complement the omitted objects, and then annotate each object with fine-grained relationships and attributes. The accomplished scene graph annotations consist of 118K objects, 83K relationships, 1M attributes, and 16K natural language actions across 43K representative video frames.

Finally, we handcraft a variety of templates to generate linguistically diverse QA pairs with both grammatical and logical guarantees. By composing the elements in the scene graphs and then filling them into proper template slots, we obtain 1.4B unbalanced and 13.4M balanced QA pairs.

3.2. Fine-grained Video Scene Graph Annotation

Representative frames. Annotating a scene graph over all video frames is impractical. Similar to [8], each of our



(a) object distribution (b) relationship distribution

(c) attribute hierarchy

Figure 2. **Statistics of the annotated video scene graphs**. We visualize the distributions of the top-15 (a) object occurrences and (b) relationship occurrences. The attributes form a hierarchical taxonomy shown in (c), where the values in the parentheses indicate the number of bottom-level attributes to be annotated. More details are provided in the supplementary material.

scene graph is annotated over a few representative frames in a video. Concretely, we use the selected frames from ANet-Entities [50] as the initialization, which cover the key semantics of all the action segments in ANet-Captions [18]. After that, we manually check and filter out those frames that hamper further annotation, *i.e.*, the frames do not contain any meaningful objects or contain too many objects from the same class. Finally, we obtain 43K frames for further annotation, which indicates an average number of 3.69 frames per video³.

Objects. ANet-Entities also provides object-level annotations for all the selected frames. Each object is annotated with a bounding box and a noun phrase (*e.g.*, "*a young woman*", "*a black jacket*"). To better organize the object annotations, we first extract nouns from the noun phrases and convert them into a set of object labels. After that, we merge the synonymous object labels (*e.g.*, "*mountain*" and "*hill*", "*saxophone*" and "*sax*"). Finally, we ask annotators to go through all the selected frames to refine the annotations, including object augmentation, label correction, and bounding box calibration. By doing the above, we obtain a total number of 118K objects of 2,072 classes over the selected frames. The top most frequent classes are shown in Figure 2a. We exclude the most frequent class "*person*" for better visualization.

Relationships. Beyond recognizing objects, predicting pairwise relationships between two objects is also important for scene understanding. Referring to the taxonomy in AGQA, we design a set of 86 relationships containing 81 contact relationships (*e.g.*, "*holding*", "*riding*", "*wearing*"),

4 spatial relationships ("near", "on", "in", "part of")⁴, and 1 temporal relationship ("identical"). Our contact relationship categories are broader than AGQA (81 vs. 16), because: (i) our videos contain both indoor and outdoor scenarios while AGOA only contains indoor ones; (ii) our relationships contain interactions between two arbitrary objects (i.e., human-object, human-human, and object-object interactions) while AGQA only contains human-object interactions. For each paired objects in one frame, annotators are asked to label at most one spatial relationship and one contact relationship, respectively. The "identical" temporal relationship indicates the objects in different frames refer to the same instance, which is used to provide indirect references of objects during question generation. Unlike other manually annotated relationships, this relationship is automatically obtained from the annotated attributes, which will be described below. The relationship occurrences follow a long-tail distribution and we illustrate the top most frequent classes in Figure 2b.

Attributes. To distinguish the fine-grained discrepancies between two objects, especially when they share the same object label, we need attribute annotations. Different from the single-label object taxonomy, the attribute taxonomy has a *multi-label* nature in that each object has multiple attributes. Moreover, the attributes for different objects are different. To address the challenges above, we handcraft a *hierarchical* attribute taxonomy by taking the characteristics of our annotated objects into consideration. As shown in Figure 2c, our attribute taxonomy includes three levels. At the top level, we categorize all the object classes into the *human* and *non-human* groups. For each group at the middle level, we design a set of representative attribute types (*e.g., "hair style"* and "*skin color"* for the *human*

 $^{^{3}}$ The number of sampled frames in our ANetQA is much lower than that of AGQA (3.69 *vs.* 24.4 on average). The motivation derives from our observation that the scene graph elements barely change within an action segment. With a limited annotation budget, we favor the annotation *density* in one frame rather than the annotation *scale* across many frames.

⁴As the viewpoints of our videos are varied, we exclude two spatial relationships ("*in front of*" and "*behind*") in AGQA to avoid ambiguity.

group, "*shape*" and "*material*" for the *non-human* group). A few attribute types like "*location*" and "*status*" are shared across the two groups. At the bottom level, we provide a set of attribute labels for each attribute type (*e.g.*, "*long hair*" and "*short hair*" for the *hair length* attribute type). For each object, annotators are asked to label the bottom-level attributes thoroughly. Due to space limitations, we only show the numbers of attributes at the bottom level in the figure. We have annotated 1M attributes over 118K objects, with an average number of 8.6 attributes per object.

As a by-product, the annotated attributes can facilitate the annotation process of the "*identical*" relationship. Specifically, if two objects in different frames have the same object label, we calculate their overlapping ratio of the annotated attributes. The pairs that surpass a confidence threshold are manually checked to ensure correctness.

To the best of our knowledge, our benchmark is the *first* attempt to provide large-scale and hierarchical attribute annotations for grounded objects in real-world videos.

Actions. In contrast to the objects, attributes, and relationships above, the action segments over specific time intervals of the video often contain much richer semantics. Using a simple label may lose the essential semantics of the action. Therefore, we use a natural language caption to describe each action segment in detail, which has been provided in ANet-Captions [18]. However, some long captions are syntactically complex and are hard to be used for question generation. To this end, we set the maximum length of a caption to 10 and filter out those captions exceeding this threshold. This results in 16K temporally-grounded captions with an average length of 8.1 words.

3.3. Compositional QA Generation

On top of the annotated spatio-temporal scene graphs, we aim to generate massive questions for diverse reasoning abilities. As shown in Table 2, we design a set of 21 question types to cover diverse reasoning skills in varying degrees of complexities. Each question type is categorized into one of the five structures (query, verify, choose, compare, and logic), which refers to the intention of the question. To fulfill the functionality of different question type, we handcraft at least one template for each question type, resulting in 119 grammatical and logical question templates. Similar to AGQA, we design a functional program for each template that traverses and composes the elements in the scene graphs, and fills them into proper template slots to produce compositional QA pairs automatically.

Compared to the question types in AGQA, our major improvements lie in that we introduce 6 extra types with respect to attributes (*i.e.*, the types starting with 'attr' in Table 2). The annotated rich attributes enable us to design up to 101 question templates (*e.g.*, "what color is ...", "what is the shape of ..."), resulting in 612.6M

type	structure	#templ.	#unbal.	#bal.
attrRelWhat [†]	query	30	169.5M	2.63M
attrWhat †	query	15	70.4M	1.43M
relWhat	query	1	33.1M	1.01M
objRelWhere	query	2	2.5M	0.55M
objRelWhat	query	2	7.1M	0.56M
objWhere	query	1	2.9M	0.43M
objWhat	query	1	0.5M	0.14M
objExist	verify	1	51.7M	1.00M
objRelExist	verify	1	98.3M	0.94M
actExist	verify	1	0.4M	0.08M
objRelWhatChoose	choose	2	347.0M	0.57M
objWhatChoose	choose	1	180.5M	0.55M
attrRelWhatChoose [†]	choose	36	149.5M	0.42M
attrWhatChoose [†]	choose	18	85.1M	0.40M
attrCompare [†]	compare	1	138.0M	2.02M
attrSame [†]	compare	1	0.09M	0.01M
actTime	compare	1	0.01M	0.01M
actLongerVerify	compare	1	0.01M	0.01M
actShorterVerify	compare	1	0.01M	0.01M
andObjRelExist	logic	1	20.2M	0.35M
xorObjRelExist	logic	1	20.2M	0.35M
overall	-	119	1.4B	13.4M

Table 2. **Statistics of the generated questions**. Each question type belongs to a certain structure and contains at least one template. More details are provided in the supplementary material. [†]: new question types that are not supported in AGQA.

unbalanced and 6.9M balanced QA pairs. Furthermore, the attribute annotations are also used to describe objects in almost all the rest templates (*e.g.*, "*what is the relationship between the* [attribute][object] and [attribute][object]?"). The introduction of attributes not only provides a more precise description of the referred object but also increases the reasoning steps of the generated questions. It is worth noting that although we can describe an object in great detail (*e.g.*, "*a walking young woman wearing green t-shirt and sunglasses*"), this would lead to a risk of combinational explosion and affect the readability of the questions. Therefore, we set the maximum number of attributes used in each question to two.

Using the above question templates, we obtain 1.4 billion QA pairs. These QA pairs are *unbalanced* and have strong language biases that models can exploit. We conduct composite balancing strategies on both the questions and answers. Following the question structure distribution in balanced AGQA, our question balancing strategy adjusts the percentages of the query/verify/choose/compare/logic questions to 50%/15%/15%/15%/5%, as shown in Figure 3a. While maintaining these percentages above, we conduct answer balancing within each question template to make sure that its answers are uniformly distributed (unbiased). In Figure 3b, we visualize the global answer distributions of the unbalanced and balanced sets in terms of the top-50 most frequent open answers (*i.e.*, the answers to the *query*



Figure 3. **Distributions before and after balancing.** (a) The question balancing is performed on question structures to adjust the percentages of the query/verify/choose/compare/logic questions to 50%/15%/15%/15%/5%. (b) The answer balancing is conducted on each question template to make its answers follow a uniform distribution. Its effect to the global answer distribution can be observed from the change in the distributions of the top 50 most frequent open answers.

structure questions). The obtained results demonstrate the effectiveness of our balancing strategies.

Our final ANetQA benchmark contains 13.4M balanced QA pairs, which consists of 10.4M train, 1.5M val, and 1.5M test samples. We compare the question and answer length distributions of ANetQA to existing Video-QA benchmarks. The results in Figure 4a show that the ANetQA questions have a wider range of lengths and are longer on average than those of all the counterparts, showing the diversity and fine granularity of our questions, respectively. Moreover, according to these challenging questions, our answer vocabulary size is much larger than that of the counterparts (see Figure 4b), which further increases the difficulty of our benchmark.

4. Experiments

This section contains comprehensive experiments and intensive analyses of ANetQA. We conduct evaluations on several state-of-the-art models and diagnose their capabilities to deal with different question structures, semantic classes, reasoning skills, and answer types, respectively. All the models are trained on the train split, validated on the val split, and evaluated on the test split. Furthermore, we also conduct a human evaluation to see the performance gap between the top-performing models and humans. Finally, we investigate the effects of different auxiliary annotations to model performance.

4.1. Experimental Setup

Compared models. We choose three state-of-the-art models for comparison, namely HCRN [19], ClipBERT [20], and All-in-one [35]. HCRN introduces a reusable conditional relation network (CRN) module and stacks multiple CRNs in depth to integrate the motion, question, and appearance features at different levels [19]. We use its



Figure 4. **Question lengths and answer vocabulary sizes.** We compare the (a) question lengths and (b) answer vocabulary sizes of our ANetQA and some typical VideoQA benchmarks like MSVD-QA [38], MSRVTT-QA [38], ActivityNet-QA [43], and AGQA [8]. Compared to the counterparts, our questions are longer and answer vocabulary size is larger, showing the fine granularity, diversity, and difficulty of our benchmark.

default settings to extract 128 appearance features and 8 motion features, respectively.

Different from HCRN, ClipBERT and All-in-one are two Transformer-based models that incorporate vision-language pretraining (VLP) on a large-scale corpus. ClipBERT is pretrained on massive image-text pairs, which enables endto-end learning by employing a sparse sampling mechanism. We adopt its official pretrained model weights as initial and then finetune the model on ANetQA using the (4×2) sampling strategy, which means 4 segments are sampled (with 2 sampled frames in each segment) at each training step. During model testing, we sample 16 frames uniformly for each video, as recommended in [20]. All-in-one is a current top-performing VideoQA model, which is the first attempt to perform end-to-end video-language pretraining using raw video and textual signals as inputs [35]. It is pretrained directly on a large-scale video-text corpus. We finetune its base model All-in-one-B on ANetOA by randomly sampling 3 frames for each video at each training step. At inference time, we also extract 3 frames uniformly and feed them to the learned model to predict the answer.

Human evaluation. We conduct an intensive human evaluation to quantify the errors and ambiguities induced during the construction of ANetQA. As the labeling costs is unaffordable to provide a thorough evaluation over all the QA pairs, we follow [8, 10] to randomly sample 4,000 QA pairs from the test set with the following two rules: (i) each video contains at least one sample, and (ii) each question type contains at least 50 samples. Each sample is assigned to five random annotators from a diverse group to answer the question and the majority vote over their predictions is regarded as the final human answer.

The human performance reach at 84.48% on the sampled test set. We take a closer look into these 15.52% inconsistent human predictions and find that they are constituted by 0.75% annotation errors, 1.95% answer ambiguities, and

taxnomy		type prior	HCRN [19]		ClipBERT [20]		All-in-one [35]		humon
			w/	w/o	w/	w/o	w/	w/o	numan
question structures	query	1.04	21.30	19.24	23.93	16.87	25.10	18.40	92.92
	compare	49.70	55.66	50.01	55.62	50.06	54.41	50.06	81.34
	choose	29.13	63.97	67.37	69.51	66.17	70.39	67.00	71.84
	verify	50.00	68.56	50.02	72.57	50.00	72.35	50.00	86.69
	logic	50.00	78.70	76.82	80.06	74.33	80.58	74.20	86.06
question semantics	object	17.74	55.99	49.55	58.69	48.22	59.81	48.99	84.26
	relationship	22.61	39.65	33.28	40.19	30.89	40.78	32.64	90.79
	attribute	14.60	35.80	34.05	39.71	32.81	40.14	33.39	82.17
	action	47.83	72.50	50.29	74.96	50.99	74.39	51.14	82.33
	object-relationship	10.48	35.17	32.38	37.66	30.03	38.42	31.32	86.47
	object-attribute	17.44	40.95	37.02	43.72	35.45	44.33	36.39	84.75
reasoning skills	duration-comparison	50.00	49.90	49.38	49.98	50.10	51.65	54.34	76.73
	exist	50.00	71.20	56.97	74.51	56.31	74.49	56.28	86.52
	sequencing	10.21	31.70	31.36	34.19	28.76	35.27	30.10	87.50
	superlative	30.32	47.46	39.78	49.55	38.83	50.14	39.60	90.14
answer types	binary	49.96	64.36	53.91	66.19	53.55	65.65	53.54	83.72
	open	6.49	29.95	29.00	33.17	26.86	34.33	28.25	84.82
overall		17.66	41.15	37.11	43.92	35.55	44.53	36.48	84.48

Table 3. A comprehensive comparison of three VideoQA methods on ANetQA. All results are evaluated on the test set. Apart from the overall accuracy, we follow [8] to report the per-type accuracies under different taxonomies. For each method, the variant trained with vision clues (w/) outperforms its blind counterpart without vision clues (w/o), implying that the language biases are well controlled.

12.82% human errors. These results imply that both of our scene graphs and generated QA pairs are of high quality. Furthermore, our benchmark contains difficult questions that even educated humans can not answer correctly. More analyses are provided in the supplementary material.

4.2. Main Results

We provide an intensive comparison of the state-of-theart methods on ANetQA In Table 3. Besides the overall accuracy, we follow [8] to report the per-type accuracies under different taxonomies, *i.e.*, question structures, question semantics, reasoning skills, and answer types. More detailed descriptions of the taxonomies and corresponding question templates are provided in the supplementary material. For each type, we provide a simple baseline, *type prior*, that uses the most frequent answer as the prediction.

From the results, we have the general observations as follows: (i) The All-in-one model pretrained on large videotext corpus achieves the overall best performance while using the least number of sampled frames. This suggests good video representations play a central role in VideoQA performance; (ii) the best performing model is still far from the human level, showing the difficulty of our benchmark and sufficient room for further improvements; and (iii) for each method, the variant trained with vision clues (w/) steadily outperforms its *blind* counterpart without any vision clues (w/o), indicating that the language biases are well controlled by our balancing strategies. The observations above are quite different from those on AGQA, where on their benchmark all models are on par with their corresponding blind counterparts. This can be explained that ANetQA has more unbalanced QA samples than AGQA, thus providing more room to perform thorough balancing strategies. Moreover, given the same model HCRN, its accuracy (especially the *open* answer type) on ANetQA is much lower than that on AGQA, verifying the fine-grained nature of our scene graphs elements.

Question structures and answer types. The *query* type questions are the most challenging ones as they have open answers. Among the rest four types which have limited answer choices⁵, the *compare* type questions report the lowest accuracy as they require more reasoning steps.

Question semantics. The attribute-oriented questions are the most difficult ones, as they require a more fine-grained understanding of video contents than the rest questions.

Reasoning skills. Similar to AGQA, each of our question is associated with one or more reasoning abilities necessary to answer the question. The questions requiring the *sequenc-ing* skills deliver the lowest accuracy as they require the temporal grounding ability. In contrast to the coarse action labels used in AGQA, our actions are depicted in natural language, which are more difficult to understand.

⁵The *compare*, *verify*, and *logic* type questions have binary answers. The *choose* type question conducts a comparison between [A] and [B], and the answer refers to one of the four choices: [A], [B], both, or none.

4.3. Effects of Auxiliary Annotations

All the comparative studies above only use the basic annotations (*i.e.*, the QA pairs) for model training. As all the QA pairs are automatically generated from scene graph annotations, it is natural to investigate whether and how auxiliary annotations facilitate model performance. We introduce two auxiliary annotations *scene graph statistics* and *oracle frames* to see their impacts on model performance, respectively. The results are provided in Table 4.

Scene graph statistics. The annotated scene graph of a given video contains all the necessary information to answer any questions on the video. Therefore, it is meaningful to investigate the impact of this information on model performance. The fine-grained characteristics of our scene graphs make it nontrivial to encode each scene graph into a feature bank like [12]. Alternatively, we introduce a simple statistical-based strategy to approximately represent the scene graph to a given video by extracting the top-Khigh-frequency (HF) words from all the questions on this video. The extracted HF words can be seamlessly used in any off-the-shelf model by concatenating them with the question words. We adopt HCRN [19] as the reference model and extract the top-40 HF words from different vocabularies (i.e., objects, relationships, attributes, and their combinations). These HF words are concatenated with the question words in both the training and testing phases.

From the results in the upper part of Table 4, we can see that adding HF objects or relationships solely do not bring further improvement over the reference model. This can be explained by the fact that relationships are strongly coupled with objects, using either of them solely can not provide sufficient scene graph information for the model to understand. Moreover, the model with HF attributes results in a distinct performance gain compared to the counterpart with HF objects. This observation verifies that our questions requires the abilities of fine-grained understanding and reasoning. Finally, exploiting all three types of HF information results in the best performance due to their complementary nature.

Oracle frames. As each question in ANetQA is generated from the scene graph elements in specific video frames, we denote these frames as the oracle frames for the question and investigate whether they can facilitate model performance. For each question, we inject the corresponding oracle frames into its sampled frames to ensure the necessary visual information to answer this question is provided. We use All-in-one [35] as the reference model since it uses few sampled frames and thus has a high probability of not covering the oracle frames. We have experimented with the oracle frames in the training, testing, and both phases, respectively. The results in the lower part of Table 4 show that injecting oracle frames in the training

	binary	open	overall
(a) scene graph statistics			
HRCN [19] (reference)	64.36	29.95	41.15
+ high-freq. objects (O)	65.81	29.29	41.18
+ high-freq. relationships (R)	63.84	29.21	40.48
+ high-freq. attributes (A)	67.67	32.21	43.75
+ high-freq. O+R+A	68.15	34.50	45.45
(b) oracle frames			
All-in-one [35] (reference)	65.65	34.33	44.53
+ training phase injection	66.54	35.18	45.40
+ testing phase injection	66.04	34.83	44.99
+ both phases injections	66.88	36.02	46.07

Table 4. **Effects of different auxiliary annotations.** (a) The scene graph statistics of a given video are represented as a set of high-frequency words extracted from all the questions of that video. (b) The oracle frames contain necessary visual information to answer a given question, which are injected in different phases.

and testing phases bring 0.87 and 0.46 point improvements over the reference model in terms of overall accuracy, respectively. Moreover, when oracle frames are applied to both the training and testing phases, the model performance is further improved due to their synergistic effects.

5. Conclusion and Future Work

In this paper, we present ANetQA, a challenging Video-QA benchmark that examines fine-grained compositional reasoning over untrimmed real-world videos. Benefiting from the fine-grained video scene graphs annotated by humans, ANetQA attains 13.4M balanced QA pairs, which is an order of magnitude larger than all previous VideoQA benchmarks. We provide comprehensive experiments and intensive analyses for state-of-the-art VideoQA methods, and the best-performing model showing that a fine-grained video understanding plays a vital role in our benchmark. Moreover, there remains a significant gap between the best model and humans, indicating the challenge of our benchmark while providing room for future improvements.

We will persistently improve our benchmark. *e.g.*, further reducing the language biases and answer ambiguities, and introducing more question types with diverse reasoning skills like scene-text understanding and causality inference. We hope that our ANetQA will serve as a cornerstone to facilitate future research in the video-language learning.

Acknowledgment. This work was supported in part by the NSFC (61836002, 62125201), in part by the Fundamental Research Funds for the Provincial Universities of Zhejiang (GK229909299001-001), in part by the NS-FC (62072147, 62020106007), and in part by the Zhejiang Provincial Natural Science Foundation of China (L-R22F020001, DT23F020007).

References

- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *ECCV*, pages 104–120. Springer, 2020. 2
- Yuhao Cui, Zhou Yu, Chunqi Wang, Zhongzhou Zhao, Ji Zhang, Meng Wang, and Jun Yu. Rosita: Enhancing vision-and-language semantic alignments via cross-and intra-modal knowledge integration. In *ACM MM*, pages 797–806, 2021.
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 2
- [4] Bernard Ghanem Fabian Caba Heilbron, Victor Escorcia and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In CVPR, pages 961–970, 2015. 1, 2, 3
- [5] Tsu-Jui Fu, Linjie Li, Zhe Gan, Kevin Lin, William Yang Wang, Lijuan Wang, and Zicheng Liu. Violet: End-toend video-language transformers with masked visual-token modeling. arXiv preprint arXiv:2111.12681, 2021. 2
- [6] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *EMNLP*, 2016. 2
- [7] Jiyang Gao, Runzhou Ge, Kan Chen, and Ram Nevatia. Motion-appearance co-memory networks for video question answering. In *CVPR*, pages 6576–6585, 2018. 2
- [8] Madeleine Grunde-McLaughlin, Ranjay Krishna, and Maneesh Agrawala. Agqa: A benchmark for compositional spatio-temporal reasoning. In *CVPR*, pages 11287–11297, 2021. 1, 2, 3, 6, 7
- [9] Deng Huang, Peihao Chen, Runhao Zeng, Qing Du, Mingkui Tan, and Chuang Gan. Location-aware graph convolutional networks for video question answering. In AAAI, pages 11021–11028, 2020. 2
- [10] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, pages 6700–6709, 2019. 6
- [11] Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *CVPR*, pages 2758–2766, 2017. 2
- [12] Jingwei Ji, Ranjay Krishna, Li Fei-Fei, and Juan Carlos Niebles. Action genome: Actions as compositions of spatiotemporal scene graphs. In *CVPR*, pages 10236–10247, 2020.
 8
- [13] Huaizu Jiang, Ishan Misra, Marcus Rohrbach, Erik Learned-Miller, and Xinlei Chen. In defense of grid features for visual question answering. In CVPR, pages 10267–10276, 2020. 2
- [14] Weike Jin, Zhou Zhao, Mao Gu, Jun Yu, Jun Xiao, and Yueting Zhuang. Multi-interaction network with object relation for video question answering. In ACM MM, pages 1193–1201, 2019. 2

- [15] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In NAACL-HLT, pages 4171–4186, 2019. 2
- [16] Aisha Urooj Khan, Amir Mazaheri, Niels Da Vitoria Lobo, and Mubarak Shah. Mmft-bert: Multimodal fusion transformer with bert encodings for visual question answering. In *EMNLP*, 2020. 2
- [17] Kyung-Min Kim, Min-Oh Heo, Seong-Ho Choi, and Byoung-Tak Zhang. Deepstory: video story qa by deep embedded memory networks. In *IJCAI*, pages 2016–2022, 2017. 3
- [18] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *CVPR*, pages 706–715, 2017. 3, 4, 5
- [19] Thao Minh Le, Vuong Le, Svetha Venkatesh, and Truyen Tran. Hierarchical conditional relation networks for video question answering. In *CVPR*, pages 9972–9981, 2020. 2, 6, 7, 8
- [20] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling. In *CVPR*, pages 7331–7341, 2021. 2, 6, 7
- [21] Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. Tvqa: Localized, compositional video question answering. In *EMNLP*, pages 9972–9981, 2018. 2, 3
- [22] Jie Lei, Licheng Yu, Tamara Berg, and Mohit Bansal. Tvqa+: Spatio-temporal grounding for video question answering. In ACL, pages 8211–8225, 2020. 3
- [23] Xiangpeng Li, Jingkuan Song, Lianli Gao, Xianglong Liu, Wenbing Huang, Xiangnan He, and Chuang Gan. Beyond rnns: Positional self-attention with co-attention for video question answering. In AAAI, pages 8658–8665, 2019. 2
- [24] Yi Li and Nuno Vasconcelos. Repair: Removing representation bias by dataset resampling. In CVPR, pages 9572–9581, 2019. 2
- [25] Fei Liu, Jing Liu, Weining Wang, and Hanqing Lu. Hair: Hierarchical visual-semantic relational reasoning for video question answering. In *ICCV*, pages 1698–1707, 2021. 2
- [26] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 10012–10022, 2021. 2
- [27] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *NeurIPS*, 32, 2019. 2
- [28] Tegan Maharaj, Nicolas Ballas, Anna Rohrbach, Aaron Courville, and Christopher Pal. A dataset and exploration of models for understanding video data through fill-in-theblank question-answering. In CVPR, pages 6884–6893, 2017. 3
- [29] Jonghwan Mun, Paul Hongsuck Seo, Ilchae Jung, and Bohyung Han. Marioqa: Answering questions by watching gameplay videos. In CVPR, pages 2867–2875, 2017. 2, 3
- [30] Seil Na, Sangho Lee, Jisung Kim, and Gunhee Kim. A readwrite memory network for movie story understanding. In *ICCV*, pages 677–685, 2017. 2

- [31] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 21(140):1–67, 2020. 2
- [32] Gunnar A Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In ECCV, pages 510–526, 2016. 2
- [33] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. Movieqa: Understanding stories in movies through questionanswering. In *CVPR*, pages 4631–4640, 2016. 2, 3
- [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 6000–6010, 2017. 2
- [35] Alex Jinpeng Wang, Yixiao Ge, Rui Yan, Yuying Ge, Xudong Lin, Guanyu Cai, Jianping Wu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. All in one: Exploring unified video-language pre-training. arXiv preprint arXiv:2203.07303, 2022. 2, 6, 7, 8
- [36] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *CVPR*, pages 9777–9786, 2021. 3
- [37] Junbin Xiao, Angela Yao, Zhiyuan Liu, Yicong Li, Wei Ji, and Tat-Seng Chua. Video as conditional graph hierarchy for multi-granular question answering. In *AAAI*, 2022. 2
- [38] 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 MM, pages 1645–1653, 2017. 2, 3, 6
- [39] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Just ask: Learning to answer questions from millions of narrated videos. In *ICCV*, pages 1686– 1697, 2021. 2, 3
- [40] Jianing Yang, Yuying Zhu, Yongxin Wang, Ruitao Yi, Amir Zadeh, and Louis-Philippe Morency. What gives the answer away? question answering bias analysis on video qa datasets. arXiv preprint arXiv:2007.03626, 2020. 2
- [41] Zekun Yang, Noa Garcia, Chenhui Chu, Mayu Otani, Yuta Nakashima, and Haruo Takemura. Bert representations for video question answering. In WACV, pages 1556–1565, 2020. 2
- [42] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B Tenenbaum. Clevrer: Collision events for video representation and reasoning. In *ICLR*, 2020. 2, 3
- [43] 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, 6
- [44] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *CVPR*, pages 6281–6290, 2019. 2
- [45] Zhou Yu, Jun Yu, Chenchao Xiang, Jianping Fan, and Dacheng Tao. Beyond bilinear: Generalized multimodal

factorized high-order pooling for visual question answering. *IEEE transactions on neural networks and learning systems*, 29(12):5947–5959, 2018. 2

- [46] Amir Zadeh, Michael Chan, Paul Pu Liang, Edmund Tong, and Louis-Philippe Morency. Social-iq: A question answering benchmark for artificial social intelligence. In *CVPR*, pages 8807–8817, 2019. 3
- [47] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. Merlot: Multimodal neural script knowledge models. *NeurIPS*, 34:23634–23651, 2021. 2
- [48] 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. 2, 3
- [49] Zhou Zhao, Jinghao Lin, Xinghua Jiang, Deng Cai, Xiaofei He, and Yueting Zhuang. Video question answering via hierarchical dual-level attention network learning. In ACM MM, pages 1050–1058, 2017. 2
- [50] Luowei Zhou, Yannis Kalantidis, Xinlei Chen, Jason J Corso, and Marcus Rohrbach. Grounded video description. In *CVPR*, pages 6578–6587, 2019. 3, 4
- [51] Linchao Zhu, Zhongwen Xu, Yi Yang, and Alexander G Hauptmann. Uncovering the temporal context for video question answering. *IJCV*, 124(3):409–421, 2017. 3