# TV-TREES: Multimodal Entailment Trees for Neuro-Symbolic Video Reasoning

**Kate Sanders**    **Nathaniel Weir**    **Benjamin Van Durme**

Johns Hopkins University

{ksande25, nweir, vandurme}@jhu.edu

## Abstract

It is challenging to perform question-answering over complex, multimodal content such as television clips. This is in part because current video-language models rely on single-modality reasoning, have lowered performance on long inputs, and lack interpetability. We propose TV-TREES, the first multimodal entailment tree generator. TV-TREES serves as an approach to video understanding that promotes interpretable joint-modality reasoning by producing trees of entailment relationships between simple premises directly entailed by the videos and higher-level conclusions. We then introduce the task of multimodal entailment tree generation to evaluate the reasoning quality of such methods. Our method's experimental results on the challenging TVQA dataset demonstrate *intepretable*, state-of-the-art zero-shot performance on full video clips, illustrating a best of both worlds contrast to black-box methods.

## 1 Introduction

Videos account for a large portion of content available and consumed online, but automated reasoning over semantically complex video-language data remains a challenging and under-explored problem. A popular task for assessing models' video understanding is narrative-centric video question-answering (VideoQA): Given a natural language question, a video clip of a movie or TV show, and a corresponding dialogue transcript, the goal is to return a correct natural language answer to the question using the video-text data.

Methods tackling this task (Yang et al., 2022; Li et al., 2020; Ko et al., 2023) frequently take the form of large, joint-modality transformer models. While these systems typically outperform smaller, domain-specific architectures, they inherently lack qualities necessary for robust and reliable video-language understanding. In addition to model performance often correlating with the length of the input video clip, analyses suggest their ability to
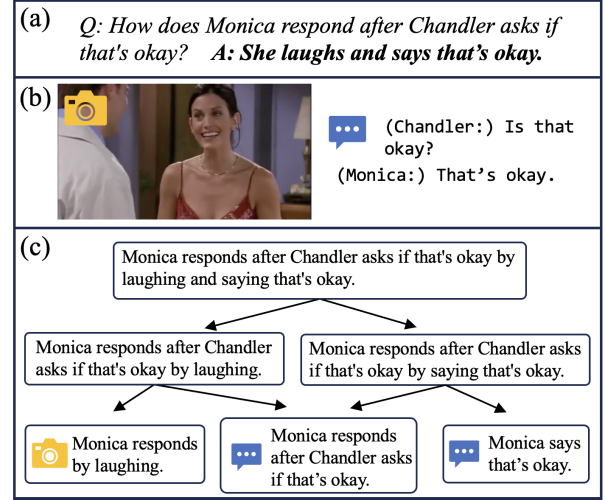


Figure 1: A (a) QA pair and (b) corresponding video clip and dialogue from the TVQA dataset (Lei et al., 2018) and (c) a multimodal entailment tree, recursively produced by our approach (top-down). Trees are created through recursively retrieving atomic evidence from the transcript and video frames and decomposing the QA pair into compositionally equivalent hypotheses until each can be directly entailed by the retrieved evidence.

perform joint visual-language reasoning is also limited, that they rely on either text or visual content but not both (Rawal et al., 2023). Better interpretability of these models could illuminate these reasoning pitfalls and allow researchers to identify and correct system issues. However, while LLMs now facilitate increasingly transparent explanation generation alongside outputs (Zhao et al., 2023), video-language models lack this ability.

Entailment trees (Dalvi et al., 2021), or trees of entailment relationships between atomic premises and higher-level conclusions, have been shown to serve well as the structural basis for text-only QA tasks by systematically and transparently modeling logical reasoning chains (Weir and Van Durme, 2023). We embrace this approach: We develop (1) the first multimodal entailment tree generator, TV-TREES (the **T**ransparent **V**ideo-**T**ext **RE**asoning

with **E**ntailment **S**ystem), and (2) the *task* of multimodal entailment tree generation to assess the reasoning ability of such systems.

In contrast to existing black-box systems, TV-TREES focuses on the manipulation of atomic "facts" retrieved from video clips to answer VideoQA questions. The approach jointly reasons over both modalities and is compatible with long video inputs, and crucially, the resulting entailment trees provide human-interpretable evidence and natural language explanations for each logical operation. Our evaluation method builds on work in informal logic and textual entailment tree generation, adapting these ideas to the multimodal domain with an emphasis on reliable evaluation.

We show that our multimodal reasoning system performs competitively on zero-shot VideoQA for the difficult TVQA dataset (Lei et al., 2018), while at the same time providing interpretable reasoning traces. Further, TV-TREES achieves state-of-the-art performance using full-length video clips as input.

In summary, our contributions are:

1. The first multimodal entailment tree generator, a fully explainable video understanding system that emphasizes logical reasoning across modalities.

2. The task of multimodal entailment tree generation and a corresponding metric for evaluating step-by-step video-text reasoning quality.

3. Results demonstrating state-of-the-art performance on zero-shot TVQA when using full clips and transcripts as input.

## 2   Related Work

### 2.1   VideoQA

QA over images makes up a large portion of multimodal question-answering work (Zou and Xie, 2020). VideoQA benchmarks constitute a smaller portion of this area (Zhong et al., 2022) and often focus on simple content and questions [CITE], but some recent videoQA datasets have targeted models' commonsense knowledge and inference ability (Lei et al., 2018; Zadeh et al., 2019). Recently, vision-and-language transformers have substantially improved performance on these videoQA tasks [CITE], and can often reason over complex content without an external knowledge base (Kim et al., 2021; Wang et al., 2021b; Salin et al., 2022).

In contrast to these video-language models, Khurana and Deshpande (2021) highlight altenrative deep learning strategies for video QA such as attention-free methods, attention-based methods, memory network methods, and hierarchical reinforced methods. Notably, Zhao et al. (2018, 2020) propose a hierarchical encoder-decoder model that uses adaptive video segmentation based on the question contents. Related works consider graph networks for video understanding (Wang et al., 2021a; Gu et al., 2021; Liu et al., 2022). While these models scale to longer videos more successfully, their performance suffers compared to transformer-based approaches.

### 2.2   Explainable Multimodal Understanding

Traditional techniques like kernel visualization and perturbation have been considered for video explainability (Hiley et al., 2019; Li et al., 2021b) alongside other approaches that consider low-level reasoning steps for simple tasks (Zhuo et al., 2019; Roy et al., 2019; Nourani et al., 2020). Some work focuses on grounded video QA, in which models are tasked with providing the visual evidence necessary for answering a question about spatial-temporal content (Xiao et al., 2023).

The approaches most similar to our work are (Chen and Kong, 2021) and (Mao et al., 2022). Chen and Kong (2021) tackle the VIOLIN video entailment dataset (Liu et al., 2020) by grounding the relevant textual entities in the video and transcript and providing a heatmap over the input as an explanation for the produced output. Our work differs in that we show exactly what data pieces contribute to the final output, explicitly model each step of the reasoning process, and don't require fine-tuning on the target dataset or domain. Mao et al. (2022) uses a chain-of-thought explanation system based on a video scene graph to answer questions about actions and objects in short video clips and GIFs. The primary difference between this and our work is the lack of dialogue and visual semantic complexity. The chain-of-thought reasoning primarily considers logical and taxonomy-centric operations over atomic-level scene graph content instead of complex inference reasoning, and the input for their proposed system only spans a few seconds.

### 2.3   Entailment Tree Generation

This paper draws inspiration from recent work on constructing natural language entailment trees to

explain reasoning. The notion starts with Dalvi et al. (2021), who introduce an expert-annotated dataset of compositional trees showing how a hypothesis follows as a logical consequence of a series of multi-premise entailment steps starting from verified support facts. They propose a series of reconstruction tasks, challenging models to reproduce expert-annotated trees given just the top-level hypothesis and some amount of gold and distractor fact leaves, and our proposed multimodal construction task is inspired by this formulation.

More recent work has introduced methods to tackle Dalvi et al.'s reconstruction task (Bostrom et al., 2022; Neves Ribeiro et al., 2022), and to use entailment trees as a basis for neuro-symbolic reasoning (Tafjord et al., 2022; Weir and Van Durme, 2023). Our work is most similar to Weir and Van Durme (2023), who introduce a QA system that reasons by searching via backward chaining for entailment trees grounded in a knowledge source. We build upon this notion, extending it to the multimodal setting and addressing the many resulting challenges.

## 2.4 Multimodal Entailment

There is a selection of work that considers entailment in images and video: (Xie et al., 2019) introduce a dataset of image-entailment pairs similar to the SNLI (Bowman et al., 2015a) corpus, and (Do et al., 2020) add natural language explanations to the pairs. More specific visual entailment tasks in this domain have been proposed as well (Thomas et al., 2022; Li et al., 2023b)., and (Suzuki et al., 2019) introduce a logic system for identifying entailment between images and captions.

Notably, Liu et al. (2020) introduce VIOLIN, a dataset of videos paired with natural language inferences that are either entailed or contradicted by the video content. Typically, standard vision-language transformers are trained for this task (Li et al., 2020; Sun et al., 2022), but more tailored approaches exist as well (Li et al., 2021a; Chen and Kong, 2021).

## 3 Multimodal Entailment Trees

We introduce the task of multimodal entailment tree generation for the VideoQA domain and the evaluation procedure.

### 3.1 Task formulation

**Input** Following Dalvi et al. (2021), as input we consider a collection of possible "evidence" and declarative form of a question-answer pair, the hypothesis $h_{(q,a)}$. Traditionally this evidence bank takes the form of a collection of natural language sentences, but in the multimodal domain, it will take the form of a video clip $V$ and corresponding dialogue transcript $D$. The video is an ordered list of $k$ images $V := \{v_i\}_{i=0}^k$, and the transcript is an ordered list of $l$ (dialogue line, timestamp) pairs $D := \{(d_i, s_i)\}_{i=0}^l$, where the timestamp maps the dialogue line to start and end frames within $V$.

**Output** We define entailment trees as structures which take the form $T := (h, e)$. $h$ is a hypothesis and $e$ is evidence, which takes the form of either a

1. *Leaf*: A (possibly empty) subset of items from $V$ or $D$.

2. *Branch*: A pair of two distinct entailment subtrees $T_1 := (h_1, e_1)$ and $T_2 := (h_2, e_2)$, where $e := (T_1, T_2)$.

Leaves with empty evidence sets are labeled as *null leaves*.

The purpose of an entailment tree is to illustrate the compositional reasoning necessary to reach a conclusion from an initial evidence bank through entailment relationships between the parent and child nodes. Therefore, in a *well-formed tree*, the evidence at any node $(h, e)$ must explicitly entail the hypothesis at that same node. For a leaf node, we posit that $e$ entails $h$ if a human would reasonably infer that $h$ is true if presented only with evidence $e \subseteq V \cup D$. For a branching node, $e$ entails $h$ if a human would reasonably infer that $h$ is true if presented with hypotheses $h_1$ and $h_2$.

**Objective** Given input $(h_{(q,a)}, V, D)$, our objective is to return a well-formed entailment tree $T$ that includes null leaves if and only if $a$ is not a correct answer to question $q$.

### 3.2 Evaluation

To serve as a secondary and distinct objective from raw VideoQA performance, we propose an evaluation method for assessing the reasoning quality of multimodal entailment trees inspired by Weir et al.'s work on scoring compositional entailments (Weir et al., 2024). Informal logic theory posits that natural language arguments may be evaluated
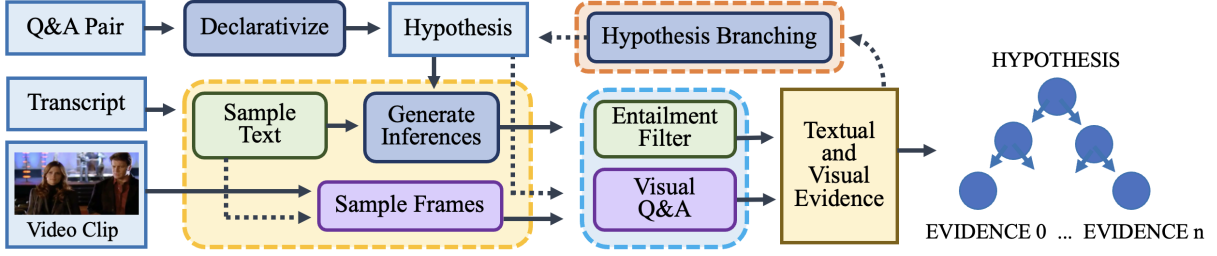
Figure 2: The multimodal proof tree generator pipeline, matching the contents of Algorithm 1. The dashed boxes divide the pipeline into the three primary modules of the system: The yellow box marks the "retrieval" module, the light blue box marks the "filter" module, and the orange box marks the "decomposition" module. With respect to individual pipeline cells, the light blue and yellow cells represent important pieces of data used or produced during the pipeline, the dark blue cells represent generative text operations, the green cells represent discriminative text operations, and the purple cells represent visual operations.

in terms of their *acceptability*, *relevance*, and *sufficiency* (Johnson and Blair, 1977). We consider each node within an entailment tree as an "argument" and consider these qualia as guidelines for comprehensive entailment tree evaluation. Below, we formulate these three qualia through an information theoretic lens to establish a set of evaluation metrics. We use the Shannon definition of information gain,

$$I(x \mid y) = -\log P(x \mid y),$$

where $P(x)$ is the probability that natural language statement $x$ is true conditioned on natural language statement(s) $y$.

**Acceptability**   Hypotheses at every node should be complete and verifiable natural language statements that are understandable to a human, and hypotheses at leaf nodes should be factually accurate statements conditioned on the world state $(V, D)$. These items may be formalized as

$$I(h) \in [0, 1] \ \ \forall h \in T \tag{1}$$
$$I(h \mid V \cup D) = 0 \ \ \forall h \in T_{\text{leaves}}. \tag{2}$$

**Relevance**   For each branching node $T_0 := (h_0, (T_1, T_2))$, hypotheses $h_1$ and $h_2$ should both be *conditionally relevant* to $h_0$, meaning that they each introduce distinct information that contributes to the compositional entailment of $h_0$. Formally, this metric is met if

$$I(h \mid h_1, h_2) < I(h \mid h_2) \ \ \forall (h, e) \in T_{\text{branches}} \tag{3}$$
$$I(h \mid h_1, h_2) < I(h \mid h_1) \ \ \forall (h, e) \in T_{\text{branches}} \tag{4}$$

**Sufficiency**   For each branching node $T_0 := (h_0, (T_1, T_2))$, hypotheses $h_1$ and $h_2$ should com-

positionally entail $h_0$, or

$$I(h_0 \mid h_1, h_2) = 0 \ \ \forall (h_0, (T_1, T_2)) \in T. \tag{5}$$

We explore practical implementations of these metrics in Section 5.

## 4   TV-TREES

In this section we introduce our proposed multimodal entailment tree generator, beginning with an overview of the framework and then individual module details. All LLM and VLM prompts are included in full in Appendix A.

---

**Algorithm 1** Tree generation, GENERATE

**Input:** Hypothesis $h$, transcript sample $D' \subseteq D$, video sample $V' \subseteq V$, current depth $k$
**Output:** Tree candidate $\hat{T} := (h, p')$
1: $F_D \leftarrow \text{RETRIEVE}(D' \mid h)$
2: $F_D' \leftarrow \text{FILTER}_D(F, h)$
3: **if** $F_D' \neq \emptyset$ **then**
4: \quad $e \leftarrow \text{BEST}_D(F_D' \mid h)$
5: **else if** $k \geq k'$ **then**
6: \quad $e \leftarrow \emptyset$
7: **else**
8: \quad $h_0, h_1 \leftarrow \text{DECOMPOSE}(h \mid T')$
9: \quad $T_0 \leftarrow \text{PROVE}(h_0, D', V', k+1)$
10: \quad $T_1 \leftarrow \text{PROVE}(h_1, D', V', k+1)$
11: \quad $e \leftarrow (T_0, T_1)$
12: **end if**
13: $F_V' \leftarrow \text{FILTER}_V(V' \mid h)$
14: **if** $\text{NULL}(e)$ and $F_V' \neq \emptyset$ **then**
15: \quad $e \leftarrow \text{BEST}_V(F_V' \mid h)$
16: **end if**
17: **return** $(h, e)$

---

4

```
Where did Beckett say they found the husband's briefcase after she said
Cynthia got rid of the gun?

(Beckett:) Actually, it's not.
(Beckett:) You got rid of the gun, of course,
(Beckett:) but we found your husband's briefcase
(Beckett:) hidden in your building's basement.
(Beckett:) There were traces of your husband's blood on it.
(Cynthia:) I gave him that briefcase
```
⬇
```
1) Beckett mentioned that they found the husband's briefcase hidden in
the building's basement.
2) There were traces of the husband's blood found on the briefcase.
3) Cynthia admitted to giving the husband the briefcase.
4) Beckett stated that Cynthia got rid of the gun.
5) The husband's briefcase was found after Cynthia got rid of the gun.
```
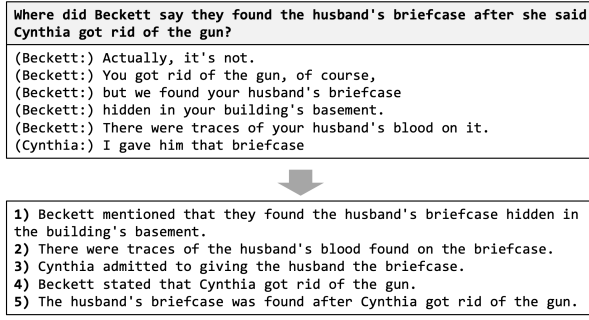
Figure 3: An example question from TVQA, corresponding dialogue excerpt sampled by TV-TREES, and set of inferences generated from these inputs by TV-TREES. The objective of inference generation is to produce a set of true natural language statements that can help prove the hypothesis.

## 4.1 System overview

TV-TREES is a recursive search algorithm that involves three primary procedures:

1. **Retrieval** Given a hypothesis and a collection of potential evidence, the system first samples relevant evidence from this collection that may sufficiently entail the current hypothesis.

2. **Filtering** The system tests whether any retrieved evidence fully entails the hypothesis. If such evidence exists and was retrieved, it is returned and the current node becomes a leaf.

3. **Decomposition** If the retrieval and filtering steps result in insufficient evidence, the system decomposes the hypothesis into two sub-hypotheses such that proving both independently is equivalent to proving the original hypothesis.

The interaction of these three parts is illustrated in Algorithm 1. Given a hypothesis $h$, transcript sample $D' \subseteq D$ and video sample $V' \subseteq V$, the system first returns evidence from the transcript relevant to $h$ (line 1) and identifies whether any of it entails $h$ (2). If such evidence was retrieved, $e$ is set to the best sample (4) and the leaf node is returned (17). Otherwise, $h$ is decomposed into sub-hypotheses $h_0$ and $h_1$ (8) and the algorithm is recursively called on these newly constructed sub-problems (9-10), treating the generated sub-proofs as explanation $e$ (11). If textual evidence cannot be found for the current node nor any of the downstream nodes (14), then the visual evidence in sample $V'$ is sampled, filtered, (13) and assigned

to $e$ where applicable (15) in the same manner as the text content.

If the maximum depth is reached during recursion, the evidence at that node is set to the empty set and the tree is incomplete.

In the following sections, we explain the implementation of the subroutines called by Algorithm 1.

## 4.2 Preprocessing

**Hypothesis Generation** The purpose of the hypothesis generation is to provide the downstream modules with a single declarative statement that contains the full semantic meaning of the original QA pair. For simplicity, this generative operation is carried out by prompting GPT-3.5 (Brown et al., 2020). We find that less robust in-context learning models like FLAN-T5 (Chung et al., 2022) are prone to omitting contextual details present in the question and not handling typos appropriately.

**Evidence Localization** Given the hypothesis, TV-TREES attempts to identify a temporal window to sample evidence from based on the dialogue. We use a cross-encoder model trained on the MS MARCO passage ranking task (Bajaj et al., 2016) to rank six-line transcript passages on their computed similarity with the generated hypothesis. We use a sliding window to calculate scores for every potential sample and return the highest scoring excerpt. If a sufficient window is identified, the vision pipeline inherits this window. If no sufficient dialogue sample is found, the system uses all video frames as the evidence bank, omitting text entirely.

## 4.3 Evidence Retrieval

Existing natural language inference (NLI) models are not well-suited for classifying entailments within highly contextual and social dialogue, which often insinuate meaning not directly stated within the text. Instead of producing an entirely new dataset for the domain of dialogue NLI, we use GPT-3.5 to generate a set of natural language inferences about the dialogue sample written in the style as data points in a dataset akin to SNLI (Bowman et al., 2015b), conditioned on a question form of the hypothesis, $q$. Presenting the question under discussion in the interrogative form significantly reduces the hallucination rate compared to passing in the original hypothesis. $q$ is also generated via GPT-3.5 taking the hypothesis $h$ as input.

5

Our system queries GPT for five inferences from a given question and passage. Then, we run these inferences through GPT to verify that they are entailed by the transcript. Examples of generated inferences are included in Figure 3.

### 4.4 Evidence Filtering

We use a cross-encoder trained on SNLI and MultiNLI to determine whether any of the retrieved evidence sufficiently entails the hypothesis. We accept any sample that achieves a logits score above a certain threshold for the "entailment" label.

Then, we apply a secondary entailment filter that ensures the inferences are accurate descriptions of the content presented in the dialogue. This is important as, while conditioning the inference generator on an interrogative form of the hypothesis mitigates hallucinations, it does not eliminate them entirely. Identifying these cases is attempted through a GPT filter that takes in the inference and the dialogue, without any hypothesis conditioning.

Finally, as the cross-encoder tends to ignore negation, which is often present in the generated inferences, we additionally pass the filtered inference-hypothesis pairs to a GPT-3.5 prompt that verifies the entailment.

The system only retains the inferences that pass through all three filters.

### 4.5 Decomposition

In the case where no atomic evidence can be retrieved from the transcript or video that immediately entails the current hypothesis, the system attempts to break it down into two sub-hypotheses that are (1) complete sentences without ambiguous pronouns or decontextualized references and (2) compositionally equivalent to the original hypothesis, i.e., proving the two sub-hypotheses as true is approximately logically equivalent to proving the original hypothesis.

We prompt GPT-3.5 to break the current hypothesis into two compositionally equivalent pieces of information, conditioned on the dialogue sample extracted in section 4.2. We instruct GPT to only return a decomposition when it is syntactically possible, to avoid recursing on sentence fragments and hypothesis repeats that the model may erroneously output if a sound decomposition cannot be found.

### 4.6 Visual Reasoning

We pass in the questions generated in Section 5.3 alongside video frames from the localized evidence

window (if applicable) sampled at 2 FPS into a vision-language model. In our experiments, we use LLaVA-7B (Liu et al., 2023). To encourage conservative classifications, in addition to asking for "yes" and "no" answers we encourage the model to respond with "not enough information" if it is unsure or the image does not provide sufficient evidence. If more than 10% of the frames in the window result in an affirmative answer from the VLM model, the visual content is considered to contain sufficient entailing evidence and the frame with the highest logits score is returned. If no frames result in an affirmative answer, no appropriate visual evidence entails the hypothesis. The LLaVA-7B prompt is included alongside the GPT prompts in Appendix A.

We also use GPT-3.5 to anonymize the question generated in section 4.1, replacing character names with common nouns such as "person". We query LLaVA-7B on each frame individually, using the anonymized question as textual input. We compare the performance of this approach to providing the original question in Appendix B, but find that the modification makes marginal difference (approximately one-point lower performance on average).

## 5 Evaluation Methodology

Traditionally, qualitative natural text evaluations have often been conducted using humans (Celikyilmaz et al., 2021), either expert annotators or crowd-sourced workers. Recently, researchers have considered whether these human evaluations could be replaced by high-performing LLMs like GPT-4 (Naismith et al., 2023). Following this line of thinking, in this section, we detail how we implement the evaluation metrics described in Section 3.2 through human annotations as well as GPT-4. We report evaluation statistics for both methods in Section 6.

### 5.1 Human Evaluations

Considering the three evaluation metrics described in Section 3.2 (acceptability, relevance, and sufficiency), we evaluate trees along these qualia through three annotation tasks. The first task provides annotators with the visual or text evidence assigned to the leaf nodes by the algorithm and ask them to assess the correctness of the leaf node hypotheses on a scale of 1-5 (**acceptability**) based on that evidence. The second task provides annotators with $(h_0, h')$ pairs from branching nodes and asks

| Method | Zero-Shot | Full Clips | Transparent | Dialogue | Vision | TVQA Acc. |
|---|---|---|---|---|---|---|
| **Fine-Tuned Methods** | | | | | | |
| STAGE | No | Yes | No | Yes | Yes | 70.5 |
| HERO | No | No | No | Yes | Yes | 74.2 |
| FrozenBiLM | No | No | No | Yes | Yes | 82.0 |
| LLaMA-VQA | No | No | No | Yes | Yes | **82.2** |
| **Zero-Shot Methods** | | | | | | |
| FrozenBiLM* | Yes | Yes | No | Yes | Yes | 26.3 |
| SeVILA | Yes | Yes | No | No | Yes | 38.2 |
| VideoChat2 | Yes | Yes | No | No | Yes | 40.6 |
| TV-TREES‡ | Yes | Yes | Yes | Yes | No | 44.9 |
| **TV-TREES** | **Yes** | **Yes** | **Yes** | **Yes** | **Yes** | **49.4** |

Table 1: Table comparing various vision-text understanding models across a set of criteria including performance on the TVQA benchmark. All zero-shot methods (Zero-Shot) take in full video clips (Full Clips), but unlike the fine-tuned approaches, none except FrozenBiLM operate over both vision and dialogue modalities. Notably, TV-TREES is the only interpretable approach. Experiment results suggest that TV-TREES and TV-TREES with text input only (TV-TREES‡) outperform existing zero-shot methods on full clips. All numbers for competing approaches are as they are reported in their respective papers except for FrozenBiLM*, which we re-run on our validation subset with full clips as input. (On ground truth clip fragments, FrozenBiLM reports 59.7% accuracy). Results suggest that a more robust visual understanding module could further improve the performance of TV-TREES, seeing the baseline results achieved by models taking in vision input only.

if the child hypothesis $h'$ is relevant to the parent $h_0$ (**relevance**). The third task provides annotators with a full hypothesis triplet $(h_0, h_1, h_2)$ from a branching node with parent $h_0$ and child premises $h_1$ and $h_2$ and asks (1) whether $h_1$ and $h_2$ each introduce distinct information (the other facet of **relevance**, we also call this **distinctness** for disambiguation purposes), and (2) if $h_0$ introduces information not provided by $h_1$ and $h_2$ together, to check for entailment (**sufficiency**). Through these tasks, annotators are also asked to indicate if any of the hypotheses or premises are malformed or otherwise uninterpretable (the other facet of **acceptability**).

Every node in a multimodal entailment tree is assigned a binary score for each assessment described above (except for the correctness checks, which are collected on a scale of 1-5). We include all task instructions and layouts in Appendix D, along with more formal descriptions of the five quantitative acceptability scores.

## 5.2 GPT Evaluations

We take the qualia outlined in Section 3.2 and write three GPT-4 prompts for (1) correct leaves in the text domain, (2) correct leaves in the vision domain, and (3) the remaining three checklist items. Correctness check prompts are modality dependent as

we use GPT-4V for vision evaluations, and separate both from the remaining checklist items as only the leaves must be evaluated for evidence-centric correctness. We use the same scoring values as in the human evaluations, and pass in twelve decompositions per prompt for the text prompts. These prompts are included in full in Appendix E.

## 5.3 Tree Scoring Paradigm

We consider the mean normalized score of the three main evaluation qualia across all nodes as the overall "composition score" for each individual tree:

$$S = \frac{a + s + 0.5(d + r)}{3}$$

where $a$ is the tree's mean leaf acceptability score, $d$ is the tree's mean distinctness score, $r$ is the tree's mean relevance score, and $s$ is the tree's mean sufficiency score.

## 6 Experiments

We evaluate TV-TREES on the TVQA dataset, comparing its performance against a text-only version of the architecture and competing zero-shot VideoQA approaches. We compare all approaches in terms of QA accuracy, and compare the entailment tree generation methods in terms of tree quality as described in Section 5.

| Trees | Acceptability | Relevance | Distinctness | Sufficiency | Score |
|---|---|---|---|---|---|
| **GPT-4 Evaluations** | | | | | |
| **Text Only** | 58.4 | 99.6 | 87.7 | 88.6 | 74.3 |
| **Multimodal** | 61.0 | 99.6 | 90.6 | 93.9 | 77.8 |
| **All** | 59.7 | 99.6 | 89.1 | 91.2 | 76.0 |
| **Human Evaluations** | | | | | |
| **Text Only** | 65.6 | 93.9 | 88.8 | 93.6 | 78.9 |
| **Multimodal** | 51.8 | 98.1 | 91.2 | 92.8 | 72.9 |
| **All** | 58.7 | 96.0 | 91.7 | 93.2 | 75.9 |

Table 2: Entailment tree quality evaluations using human and LLM evaluators. For the human annotations, acceptability corresponds to Task 1, relevance to Task 2, and distinctness and sufficiency to Task 3. These metrics are explicitly labeled in the GPT-4 prompts for evaluation. This table reports mean scores aggregated per tree for each category. In addition to metric scores, we report composition score as defined in Section 5.3. We partition results by modality: We report scores for trees using text content only, trees that use visual evidence, and both groups combined. As shown, tree scores largely suffer due to the correctness of the leaf nodes, which is unsurprising given the difficulty of extracting high-level inferences from social dialogue and often ambiguous video screenshots.

| Method | Acc. | Comp. Acc. | Comp. % |
|---|---|---|---|
| Vision | 32.4 | 51.9 | 19.7 |
| Dialogue | 44.9 | 53.3 | 51.5 |
| Both | 49.4 | 53.0 | 69.5 |

Table 3: Ablation experiment results comparing performance on TVQA when using only dialogue evidence, only visual evidence, and both modalities as evidence. We report the overall accuracy, the accuracy of the system on questions where at least one proof was complete, and the percentage of questions on which at least one proof was complete.

## 6.1 Setup

We instantiate TV-TREES as it is described in Section 4, setting the maximum recursion depth to $k = 2$, or allowing trees with up to 3 levels. Our experiments focus on the multiple choice VideoQA domain, and so we consider a question's correct answer to be the answer that results in a complete tree. In the case that the system does not successfully complete any tree for the five answer candidates, we consider the answer candidate with the "most complete" tree to be the correct answer, breaking ties by the average entailment score at each node. When complete trees are generated for multiple answers, we break ties in the same way.

## 6.2 Evaluation on TVQA

**Data** We evaluate our system on 3,000 multiple choice questions from the validation set of TVQA (Lei et al., 2018). TVQA is a VideoQA benchmark that includes multiple choice questions about the dialogue and visual content of video clips taken from six TV shows. The clips are approximately 60-90 seconds long and contain around 30 lines of dialogue each. An example TVQA question is shown in Figure 1.

**Models** In the zero-shot setting, in addition to TV-TREES, we consider zero-shot approaches FrozenBiLM (Yang et al., 2022), SeVILA (Yu et al., 2023), and VideoChat2 (Li et al., 2023a). We also include performance reported by other systems (not zero-shot) for context: STAGE (Lei et al., 2019), HERO (Li et al., 2020), FrozenBiLM (fine-tuned) (Yang et al., 2022), and LLaMA-VQA (Ko et al., 2023).

**Ablations** Existing work notes that both existing multimodal models are biased toward the text modality, often relying on text data for reasoning even for video-centric questions. In line with this theme, we evaluate our system's performance conditioned on input modality on a subset of the TVQA validation set. We first evaluate the system when it is only provided with dialogue transcripts from the clip and then when it is only provided with video frames from the clip.

**Results** We report overall accuracy alongside qualitative comparisons between the approaches in Table 4. As shown in the table, TV-TREES outperforms existing zero-shot methods when using full clips, but still shows significant room for future improvements. Notably, the text-only model

outperforms joint-modality methods, and the joint modality model only improves performance modestly, suggesting that the language modules of TV-TREES are more robust and performance could be further increased through improvements to the vision pipeline. This is further shown in the ablation experiment results in Table 3, which suggests that vision evidence alone allows TV-TREES to complete trees for only 19.7% of the questions compared to 51.5% and 69.5% for text-only and joint-modality models, respectively.

## 6.3 Proof Scoring

**Setup** We randomly sample 600 completed entailment trees generated by TV-TREES on the TVQA validation split, split evenly between text-only and multimodal trees and split evenly among tree complexity (ranging from one to seven tree nodes).

We evaluate these sampled trees using the automatic GPT4 approach as described in Section 5.2. We then sample 200 proofs from this set (evenly distributed across modalities and complexity) and we annotate this set with human annotators from Amazon Mechanical Turk as described in Section 5.1. For human annotations, we identify careful annotators through a preliminary pilot task where each annotator's work is scored by hand, and only high-scoring annotators are invited to annotate the full proofs. More information regarding these crowd-sourced annotations are included in Appendix C.

For scoring acceptability, we provide the scorer with the localized dialogue retrieved by the cross encoder model described in Section 4.2, and the video frames that achieved the highest logits scores during VQA inference, depending on the modality. We report results in Table 2.

**Results** Generally, there is a close alignment between the GPT-4 and human scores. While the overall average score assigned to the trees is within a .1 point difference between the two approaches, GPT-4 tended to score the text-only trees more harshly than humans, and the multimodal trees more leniently. This is shown primarily in the resulting acceptability scores, and more moderately in the sufficiency scores. GPT-4 rated relevance scores more leniently for both modalities, which may stem from differences in human interpretations of the task instructions. In contrast, distinctness scores are almost identical between the two methods.

We find that, unsurprisingly, the majority of errors in the produced trees stem from acceptability issues. According to human evaluations, the visual module produces lower quality inferences than the textual modules do. This is not surprising, as we are able to include additional entailment filters for the textual reasoning steps to remove lower quality predictions before constructing the final entailment trees, whereas we do not have similar methods in place for visual inference. Based on these results, introducing stronger entailment classifiers for both domains may significantly improve performance on tree evaluation as well as on general VideoQA.

## 7 Conclusion

We introduce the first neuro-symbolic entailment tree generator for multimodal content to improve robustness, reliability, interpretability, and scalability of video-language understanding systems. We focus on the application of narrative-driven VideoQA, and show that our approach achieves state-of-the-art results on the zero-shot TVQA benchmark with full video clips. We also propose the *task* of multimodal entailment tree generation for the assessment of generated tree reasoning quality, establishing an information-theoretic evaluation method grounded in informal logic theory. Experimental results suggest that such interpretable, neuro-symbolic approaches to video understanding are a strong alternative to existing methods and present exciting directions for future research.

## 8 Limitations

We introduce an initial exploration into the task of multimodal entailment tree generation for video understanding, and so, there are inherent limitations that we hope to correct in future work. Most notably, our vision module underperforms compared to some systems - in future work, we hope to improve upon the existing end-to-end architecture as well as explore more compositional approaches. Furthermore, while we consider six lines of dialogue at a time to ensure sufficient context for textual inference, we do not do the same for visual analysis (instead working with only one frame at a time). Extending the immediate context for visual inference would likely improve performance as well. Finally, it is important to consider the domain that our system is used in, as model performance may vary in domains with limited dialogue, etc. We hope that this work inspires future research in this domain to improve upon our proposed pipeline.

# References

Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. 2016. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268*.

Kaj Bostrom, Zayne Sprague, Swarat Chaudhuri, and Greg Durrett. 2022. Natural language deduction through search over statement compositions. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4871–4883, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015a. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015b. A large annotated corpus for learning natural language inference. *CoRR*, abs/1508.05326.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. 2021. Evaluation of text generation: A survey.

Junwen Chen and Yu Kong. 2021. Explainable video entailment with grounded visual evidence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.

Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan Xie, Hannah Smith, Leighanna Pipatanangkura, and Peter Clark. 2021. Explaining answers with entailment trees. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7358–7370, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. 2020. e-snli-ve: Corrected visual-textual entailment with natural language explanations. *arXiv preprint arXiv:2004.03744*.

Mao Gu, Zhou Zhao, Weike Jin, Richang Hong, and Fei Wu. 2021. Graph-based multi-interaction network for video question answering. *IEEE Transactions on Image Processing*, 30:2758–2770.

Liam Hiley, Alun Preece, and Yulia Hicks. 2019. Explainable deep learning for video recognition tasks: A framework & recommendations. *arXiv preprint arXiv:1909.05667*.

Ralph H. Johnson and J. Anthony Blair. 1977. Logical self-defense.

Khushboo Khurana and Umesh Deshpande. 2021. Video question-answering techniques, benchmark datasets and evaluation metrics leveraging video captioning: a comprehensive survey. *IEEE Access*, 9:43799–43823.

Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning*, pages 5583–5594. PMLR.

Dohwan Ko, Ji Soo Lee, Wooyoung Kang, Byungseok Roh, and Hyunwoo J Kim. 2023. Large language models are temporal and causal reasoners for video question answering. *arXiv preprint arXiv:2310.15747*.

Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. 2018. Tvqa: Localized, compositional video question answering. *arXiv preprint arXiv:1809.01696*.

Jie Lei, Licheng Yu, Tamara L Berg, and Mohit Bansal. 2019. Tvqa+: Spatio-temporal grounding for video question answering. *arXiv preprint arXiv:1904.11574*.

Juncheng Li, Siliang Tang, Linchao Zhu, Haochen Shi, Xuanwen Huang, Fei Wu, Yi Yang, and Yueting Zhuang. 2021a. Adaptive hierarchical graph reasoning with semantic coherence for video-and-language inference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1867–1877.

Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. 2023a. Mvbench: A comprehensive multi-modal video understanding benchmark. *arXiv preprint arXiv:2311.17005*.

Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. 2020. Hero: Hierarchical encoder for video+ language omni-representation pre-training. *arXiv preprint arXiv:2005.00200*.

Nan Li, Pijian Li, Dongsheng Xu, Wenye Zhao, Yi Cai, and Qingbao Huang. 2023b. Scene-text oriented visual entailment: Task, dataset and solution. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5562–5571.

Zhenqiang Li, Weimin Wang, Zuoyue Li, Yifei Huang, and Yoichi Sato. 2021b. Towards visually explaining video understanding networks with perturbation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1120–1129.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning.

Jingzhou Liu, Wenhu Chen, Yu Cheng, Zhe Gan, Licheng Yu, Yiming Yang, and Jingjing Liu. 2020. Violin: A large-scale dataset for video-and-language inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10900–10910.

Yun Liu, Xiaoming Zhang, Feiran Huang, Bo Zhang, and Zhoujun Li. 2022. Cross-attentional spatio-temporal semantic graph networks for video question answering. *IEEE Transactions on Image Processing*, 31:1684–1696.

Jianguo Mao, Wenbin Jiang, Xiangdong Wang, Zhifan Feng, Yajuan Lyu, Hong Liu, and Yong Zhu. 2022. Dynamic multistep reasoning based on video scene graph for video question answering. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3894–3904.

Ben Naismith, Phoebe Mulcaire, and Jill Burstein. 2023. Automated evaluation of written discourse coherence using gpt-4. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 394–403.

Danilo Neves Ribeiro, Shen Wang, Xiaofei Ma, Rui Dong, Xiaokai Wei, Henghui Zhu, Xinchi Chen, Peng Xu, Zhiheng Huang, Andrew Arnold, and Dan Roth. 2022. Entailment tree explanations via iterative retrieval-generation reasoner. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 465–475, Seattle, United States. Association for Computational Linguistics.

Mahsan Nourani, Chiradeep Roy, Tahrima Rahman, Eric D Ragan, Nicholas Ruozzi, and Vibhav Gogate. 2020. Don't explain without verifying veracity: an evaluation of explainable ai with video activity recognition. *arXiv preprint arXiv:2005.02335*.

Ishaan Singh Rawal, Shantanu Jaiswal, Basura Fernando, and Cheston Tan. 2023. Revealing the illusion of joint multimodal understanding in videoqa models. *arXiv preprint arXiv:2306.08889*.

Chiradeep Roy, Mahesh Shanbhag, Mahsan Nourani, Tahrima Rahman, Samia Kabir, Vibhav Gogate, Nicholas Ruozzi, and Eric D Ragan. 2019. Explainable activity recognition in videos. In *IUI Workshops*, volume 2.

Emmanuelle Salin, Badreddine Farah, Stéphane Ayache, and Benoit Favre. 2022. Are vision-language transformers learning multimodal representations? a probing perspective. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11248–11257.

Yuchong Sun, Hongwei Xue, Ruihua Song, Bei Liu, Huan Yang, and Jianlong Fu. 2022. Long-form video-language pre-training with multimodal temporal contrastive learning. *Advances in neural information processing systems*, 35:38032–38045.

Riko Suzuki, Hitomi Yanaka, Masashi Yoshikawa, Koji Mineshima, and Daisuke Bekki. 2019. Multimodal logical inference system for visual-textual entailment. *arXiv preprint arXiv:1906.03952*.

Oyvind Tafjord, Bhavana Dalvi Mishra, and Peter Clark. 2022. Entailer: Answering questions with faithful and truthful chains of reasoning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2078–2093, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Christopher Thomas, Yipeng Zhang, and Shih-Fu Chang. 2022. Fine-grained visual entailment. In *European Conference on Computer Vision*, pages 398–416. Springer.

Jianyu Wang, Bing-Kun Bao, and Changsheng Xu. 2021a. Dualvgr: A dual-visual graph reasoning unit for video question answering. *IEEE Transactions on Multimedia*, 24:3369–3380.

Zekun Wang, Wenhui Wang, Haichao Zhu, Ming Liu, Bing Qin, and Furu Wei. 2021b. Distilled dual-encoder model for vision-language understanding. *arXiv preprint arXiv:2112.08723*.

Nathaniel Weir, Kate Sanders, Orion Weller, Shreya Sharma, Dongwei Jiang, Zhengping Zhang, Bhavana Dalvi Mishra, Oyvind Tafjord, Peter Jansen, Peter Clark, et al. 2024. Enhancing systematic decompositional natural language inference using informal logic. *arXiv preprint arXiv:2402.14798*.

Nathaniel Weir and Benjamin Van Durme. 2023. Dynamic generation of grounded logical explanations in a neuro-symbolic expert system.

Junbin Xiao, Angela Yao, Yicong Li, and Tat Seng Chua. 2023. Can i trust your answer? visually grounded video question answering. *arXiv preprint arXiv:2309.01327*.

Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. 2019. Visual entailment: A novel task for fine-grained image understanding. *arXiv preprint arXiv:1901.06706*.

Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. Zero-shot video question answering via frozen bidirectional language models. *Advances in Neural Information Processing Systems*, 35:124–141.

Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. 2023. Self-chained image-language model for video localization and question answering. *arXiv preprint arXiv:2305.06988*.

Amir Zadeh, Michael Chan, Paul Pu Liang, Edmund Tong, and Louis-Philippe Morency. 2019. Social-iq: A question answering benchmark for artificial social intelligence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8807–8817.

Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2023. Explainability for large language models: A survey. *ACM Transactions on Intelligent Systems and Technology*.

Zhou Zhao, Shuwen Xiao, Zehan Song, Chujie Lu, Jun Xiao, and Yueting Zhuang. 2020. Open-ended video question answering via multi-modal conditional adversarial networks. *IEEE Transactions on Image Processing*, 29:3859–3870.

Zhou Zhao, Zhu Zhang, Shuwen Xiao, Zhou Yu, Jun Yu, Deng Cai, Fei Wu, and Yueting Zhuang. 2018. Open-ended long-form video question answering via adaptive hierarchical reinforced networks. In *IJCAI*, volume 2, page 8.

Yaoyao Zhong, Junbin Xiao, Wei Ji, Yicong Li, Weihong Deng, and Tat-Seng Chua. 2022. Video question answering: Datasets, algorithms and challenges. *arXiv preprint arXiv:2203.01225*.

Tao Zhuo, Zhiyong Cheng, Peng Zhang, Yongkang Wong, and Mohan Kankanhalli. 2019. Explainable video action reasoning via prior knowledge and state transitions. In *Proceedings of the 27th acm international conference on multimedia*, pages 521–529.

Yeyun Zou and Qiyu Xie. 2020. A survey on vqa: Datasets and approaches. In *2020 2nd International Conference on Information Technology and Computer Application (ITCA)*, pages 289–297. IEEE.

## A  TV-TREES LLM Prompts

We provide the LLM and VLM prompts used in the TV-TREES pipeline in Figures 9-16.

## B  Visual Prompt Anonymization Experiments

We consider an additional component to the TV-TREES system outlined in Section 4 that anonymizes any references to characters passed in to the visual entailment module. We pass any questions that will be used for visual QA prompts through a GPT filter that replaces any character names with common nouns and pronouns like "the man", "they", and "the doctor". We report results below, comparing this alternate system to the competing methods and the standard TV-TREES method. We find that the anonymization paradigm results in a TVQA accuracy score of 48.1% compared to the standard system's 49.4%. We provide the anonymization GPT prompt in Figure 13 and a results table for comparison (Table 4).

## C  Amazon Mechanical Turk Details

We evaluate generated tree quality through crowdsourced workers on Amazon Mechanical Turk with three main annotation tasks. We identify a separate group of quality annotators for each task by (1) setting the qualifications for the task to workers located within the United States with a HIT acceptance rate of 98% and over 1000 completed HITS, and (2) running a pilot task with carefully selected questions to identify annotators who answer the preselected questions with high accuracy.

   We estimate time completion for each version of the task uploaded to Mechanical Turk and set the payment values to an estimated $15 per hour. No identifiable information of any annotators is present in this paper or in any artifacts we will release.

## D  Human Tree Evaluation Tasks

Below, we include screenshots depicting the instructions and format of each task provided to annotators. We also include a table detailing the descriptions provided to annotators for each of the five acceptability scores (Table 5).

**Acceptability:** See Figures 4 and 5.

**Relevance:** See Figure 6.

**Sufficiency:** See Figures 7 and 8.

## E  GPT-4 Evaluation Prompts

Prompts for GPT-4 evaluations are shown in Figures 17 - 19. Figure 17 shows the primary decomposition evaluation prompt, which accounts for relevancy, distinctness, and sufficiency. Figure 18 shows the textual acceptability for dialogue prompt, and Figure 19 shows the visual acceptability for screenshots prompt, which was passed to GPT-4V.

Your task is to rate how correct statements about TV show clips are based on a small amount of evidence taken from the show.

You will be given a statement about a TV show clip and dialogue snippet from the TV clip. **This version of the task does not use a screenshot from the show - base your answers on the dialogue only.** Drag the slider to indicate whether or not the statement is true based only on the provided transcript. The sliding scale has five options to allow for cases where you can't determine the statement's factuality, or otherwise aren't completely sure.

Check "The statement doesn't make sense." if the statement is not a complete sentence or is otherwise impossible to assess the correctness of due to grammatical issues, etc.

**Thank you for participating in the task. Quality of submitted tasks will be carefully monitored, and bot outputs will be rejected.**

---

**Dialogue:**

I could act like a writer. Here...
But see? Nada.
(Joey:)I don't have the discipline
that it takes. I can't do it.
(Ross:)I'll help you.
(Ross:)Yeah. I'll make up a schedule
and make sure you stick to it.
(Ross:)Plus, it'll give me something to do.

*Joey says he doesn't have the discipline to be a writer.*

Is the sentence above an accurate statement about the TV script?

False ●————————————————— True
–
☐ The statement doesn't make sense.

---

Figure 4: AMT acceptability task instructions and example for premises with textual evidence.

**If the statement mentions a character whose face you do not know, you can click on the link directly below the screenshot to see a list of faces and names for that show. From there, you can press CTRL+F and search for the name you aren't sure about.**

Your task is to rate how correct statements about TV show clips are based on a small amount of evidence taken from the show.

You will be given a statement about a TV show clip and a corresponding screenshot and dialogue snippet from the TV clip. Drag the slider to indicate whether or not the statement is true based only on the provided screenshot and transcript. The sliding scale has five options to allow for cases where you can't determine the statement's factuality, or otherwise aren't completely sure.

**Thank you for participating in the task. Quality of submitted tasks will be carefully monitored, and bot outputs will be rejected.**

---



Link to the character faces & names

*The color of the button by the door when they walk through is orange.*

Is the sentence above an accurate statement about the TV screenshot and/or script?

False ●————————————————— True
–
☐ The statement doesn't make sense.

---

Figure 5: AMT acceptability task instructions and example for premises with visual evidence.

Your task is to label whether two sentences are related to each other, and whether they contradict each other.

In this task, facts are "related" if they both refer to one or more of the same people, places, or things.
Facts are "contradicting" if they say two things that cannot be true or are unlikely to be true at the same time.

**Thank you for participating in the task. Quality of submitted tasks will be carefully monitored, and bot outputs will be rejected.**

---

Fact A: *House said he wanted to get rid of the bad cells in Abigail's body by using severe chemo when talking about treatments with Abigail and Maddy.*

Fact B: *House was talking about treatments with Abigail and Maddy.*

Are Fact A and Fact B related to each other?

○ Related statements.
○ Unrelated statements.

Are Fact A and Fact B contradicting each other?

○ Contradicting statements.
○ No contradiction.

---

Figure 6: AMT relevance task instructions and example.

In this task, you will be presented with three facts, sequentially. After receiving Fact A and Fact B, you will mark whether the two facts present different information or not. For example,

"Turtles are animals" and "Turtles are reptiles"
contain different information, and

"Turtles have webbed feet and are animals" and "Turtles have webbed feet and are reptiles"
still contain different information, but

"Turtles are a type of animal" and "Turtles are animals"
contain the same information.

Then, you will receive Fact C. Considering Fact A and Fact B together, you should identify whether Fact C introduces new information not included in A and B. For example, if Fact A is "Turtles have webbed feet and are animals" and Fact B is "Turtles are reptiles", then

"Turtles are webbed-footed reptiles" does not introduce new information, but
"Turtles are webbed-footed omnivores" does.

Assume that each fact refers to the same point in time - for example, if Fact A is "Joey throws the ball" and Fact B is "Joey says 'Catch!'", then the fact "Joey says 'catch' when he throws the ball" does not provide new information.

**Check "Fact C doesn't make sense." if the statement is not a complete sentence or is otherwise impossible to assess due to grammatical issues, etc.**

**Thank you for participating in the task. Quality of submitted tasks will be carefully monitored, and bot outputs will be rejected.**

Figure 7: AMT sufficiency task instructions.

---

Here is Fact A.

Fact A: ***Castle said he was erasing his memory.***

---

Here is Fact B.

Fact B: ***The video was playing.***

Does Fact B provide any information not included in Fact A?

○ **Yes, Fact B includes new information not mentioned by Fact A.**
○ **No, Fact B provides no new information.**

---

Here is Fact C.

Fact C: ***Castle said he was erasing his memory when the video was playing.***

When considering Fact A and Fact B together, does Fact C provide any information not already covered by A and B together?

○ **Yes, Fact C includes new information not mentioned by Fact A and B.**
○ **No, Fact C does not provide any new information.**
☐ Fact C doesn't make sense.

Figure 8: AMT sufficiency task example.

| Method | FrozenBiLM | SeVILA | VideoChat2 | TV-TREES[‡] | TV-TREES | TV-TREES* |
|---|---|---|---|---|---|---|
| TVQA Acc. | 26.3 | 38.2 | 40.6 | 44.9 | 49.4 | 48.1 |

Table 4: Table contextualizing the anonymized VQA inputs ablation experiment (TV-TREES*) by comparing it to the other zero-shot TVQA results.

| Score | Description |
|---|---|
| 1 | Sentence is contradicted by the screenshot or dialogue. |
| 2 | Sentence is **unlikely** to be true based on the screenshot or dialogue. |
| 3 | Sentence is purely ambiguous given the screenshot or dialogue. |
| 4 | Sentence is **likely** to be true based on the screenshot or dialogue. |
| 5 | Sentence is directly suggested or shown by the screenshot or dialogue. |

Table 5: Descriptions for each acceptability score provided to annotators as part of the sliding bar functionality in the task.

**Hypothesis Generation Prompt**

```
Convert each of the answer options for the following questions into GRAMMATICAL
ANSWER SENTENCES. Make sure that they are FULL and COMPLETE sentences, not just
words.  They should be sentences that you can "prove" by reasoning about the
situation.  Proving the sentence should amount to choosing choosing that answer
option over the other one(s).

## Input
QUESTION:
{ICL Q Examples}

## Output
{ICL A Examples}


## Input
QUESTION:
{Questions}

## Output
```

Figure 9: Example prompt for generating hypotheses from QA pairs as described in Section 4.2.

**Hypothesis-To-Question Generation Prompt**

```
Rewrite the following statement into a "yes" or "no" question, and nothing else.

STATEMENT: "{Statement}"
QUESTION:
```

Figure 10: Example prompt for generating interrogative forms of hypotheses for conditioning inference generation and VQA as described in Section 4.3.

**Hypothesis Decomposition Prompt**

```
You are a writing system that values clarity above all else.  You NEVER uses
pronouns like "he", "they", or "it" to ensure that readers can understand your
sentences in isolation without additional context.

Your task is to break down the following statement into two, simpler sentences.

STATEMENT: "Lauren closed the door after discussing the party with Kelly."

DECOMPOSITION (USING NO PRONOUNS, INCLUDING "THEY" OR "HE" OR "SHE"):
(1) "Lauren closed the door."
(2) "Lauren discussed the party with Kelly."

STATEMENT: "Jason asked about the brown briefcase because he was concerned that it
had been misplaced or stolen."

DECOMPOSITION (USING NO PRONOUNS, INCLUDING "THEY" OR "HE" OR "SHE"):
(1) "Jason asked about the brown briefcase."
(2) "Jason was concerned that the brown briefcase had been misplaced or stolen."

STATEMENT: "{Statement}"

DECOMPOSITION (USING NO PRONOUNS, INCLUDING "THEY" OR "HE" OR "SHE"):
```

Figure 11: Example prompt for decomposing a hypothesis into two distinct premises as described in Section 4.5.

**Inference Generation Prompt**

```
You are a fact-checking expert that uses evidence to answer questions about a TV
show.

For the following question and scene dialogue, write a set of five independent
inferences entailed by some part of the scene.  The inferences should resemble
short, factual statements about the scene and should help to answer the question
using component reasoning steps.

Write your facts in JSON format, i.e.  {"1":  "<answer here>", "2":  "<answer
here>", ...} and nothing else.

QUESTION: "Why does Howard say theyŕe late after walking in?"

SCENE:
{Dialogue}

INFERENCES (5 total):
```

Figure 12: Example prompt for generating inferences from dialogue samples given an underlying question as described in Section 4.3.

**Premise-Dialogue Entailment Verification Filtering Prompt**

```
You are an expert social reasoning system that understands the implied meanings
of complex conversations between TV show characters.  Given social inferences made
by other AI systems about transcripts, you score them on whether they are CORRECT or
NOT SUPPORTED by the transcript.

Given the following TV show transcript, write whether each of the following
statements about the TV show are CORRECT or NOT SUPPORTED. A statement is CORRECT
if an average human would agree that it is most likely true based on the transcript,
and is NOT SUPPORTED otherwise.

Write your facts in JSON format, i.e.  {"1":  <"answer here">, "2":  <"answer
here">, ...} and nothing else.

TRANSCRIPT:
{Dialogue}

STATEMENTS:
{Inferences}

OUTPUT:
```

Figure 13: Example prompt for filtering premises based on dialogue entailment as described in Section 4.3.

**Question Anonymization Prompt**

```
Anonymize the following questions by replacing all the characters' names replaced
with "the man", "the woman", "the person", or "the people". Your output should be formatted
as a serialized JSON list, i.e.  {"1":  <answer here>, "2":  <answer here>}, ..., and
nothing else.

SENTENCES:
{Questions}

QUESTIONS:
```

Figure 14: Example prompt for generating anonymized versions of interrogative versions of hypotheses as described in Appendix B.

**Premise-Hypothesis Entailment Verification Filtering Prompt**

```
You are a logical reasoning system that determines whether individual facts are
enough to prove a hypothesis statement.

For each of the following independent facts, answer "YES" if the fact cannot be
true without the hypothesis also being true, and "NO" if the hypothesis can be false
even if the fact is true.  Always answer "NO" if the hypothesis is not a complete
sentence (for example "is sitting.".  Write your answers in JSON format, i.e.  {"1":
"<fact 1 answer here>", "2":  "<fact 2 answer here>", ...} and nothing else.

HYPOTHESIS: {Hypothesis}

FACTS:
{Inferences}

OUTPUT:
```

Figure 15: Example prompt for filtering premises based on hypothesis entailment as described in Section 4.4.

**Visual QA Prompt**

```
From this image, can you answer the question {Question}?  If so, answer the
question, otherwise, answer "NOT ENOUGH INFO".
```

Figure 16: Prompt template for soliciting VQA outputs from the LLaVA-7B model as described in Section 4.6.

**GPT-4 Relevance, Distinctness, and Sufficiency Evaluation**

```
You are a reasoning system that searches for proofs of a hypothesis about a video
clip by recursively decomposing it into simpler premises.

Given a hypothesis, you identify entries in a list of possible two-premise
decompositions of the hypothesis that are "well-formed": Proving the premises
of a well-formed decomposition would amount to proving the hypothesis through
compositional entailment.

You assess decompositions using three metrics: Premise relevancy, premise
distinctness, and decomposition sufficiency. Each decomposition should receive
two relevancy and distinctness scores, one for each premise, but only one single
sufficiency score.

RELEVANCY: Relevancy measures whether a premise contributes information pertaining
to the hypothesis. This is measured on a binary scale. Simply, if the premise
mentions an entity or idea also mentioned by the hypothesis, the relevancy score is
1. Otherwise, it is 0.

DISTINCTNESS: Distinctness measures whether a premise introduces new information not
already entailed by the other premise in the decomposition. This is measured on a
binary scale. If the premise only introduces information already entailed by the
other premise in the decomposition, the distinctness score is 0. Otherwise, it is 1.
If both premises are the same, both receive a score of 0.

SUFFICIENCY: Sufficiency measures whether the two premises cover all the information
introduced by the hypothesis. This is also measured on a binary scale. If, when
considering both premises, the hypothesis introduces new information not covered by
the decompositional premises, the sufficiency score is 0. If the hypothesis does
not introduce new information, the sufficiency score is 1.

For the following decompositions, score each decomposition's relevancy and
sufficiency. Decompositions will be presented in the form "(<decomposition number>)
H: <hypothesis> & P1: <decomp premise 1> & P2: <decomp premise 2>". Your answer
should be a list of entries taking the form "(<decomposition number>) RELEVANCY:
(<premise 1 score>, <premise 2 score>), DISTINCTNESS: ((<premise 1 score>, <premise
2 score>), SUFFICIENCY: (<overall score>)".

DECOMPOSITIONS:
{Decompositions}

JUDGEMENTS (one line per decomposition):
```

Figure 17: GPT-4 prompt for scoring the relevance, distinctness, and sufficiency of decompositions in an entailment tree.

**GPT-4 Textual Acceptability Evaluation**

```
Based on the dialogue from the TV show, how likely is it that the statements below
are true? Score the likelihood of each statement on a 1-5 scale, where 1 indicates
the dialogue contradicts the statement, 2 indicates the statement is unlikely to be
true given the dialogue, 3 indicates the statement is ambiguous given the dialogue,
4 indicates the statement is likely to be true given the dialogue, and 5 indicates
that the statement must be true given the dialogue. Write your numerical scores in
the same order as the listed statements, separated by commas, and nothing else.

Dialogue:
{Dialogue}

Statements:
{Statements}
```

Figure 18: GPT-4 prompt for scoring the acceptability of entailment tree leaf nodes that cite textual evidence.

**GPT-4V Visual Acceptability Evaluation**

```
Based on the screenshot from the TV show, how likely is it that the statement below
is true?  Score the likelihood on a 1-5 scale, where 1 indicates the screenshot
contradicts the statement, 2 indicates the statement is unlikely to be true given
the screenshot, 3 indicates the statement is ambiguous given the screenshot, 4
indicates the statement is likely to be true given the screenshot, and 5 indicates
that the statement must be true given the screenshot.  Write your numerical score
and nothing else.

Statement:  {Statement}
```

Figure 19: GPT-4V prompt for scoring the acceptability of entailment tree leaf nodes that cite visual evidence. The top-scoring video frame is passed in alongside the prompt.