

UC Merced

Proceedings of the Annual Meeting of the Cognitive Science Society

Title

CogME: A Cognition-Inspired Multi-Dimensional Evaluation Metric for Story Understanding

Permalink

<https://escholarship.org/uc/item/8p3137gd>

Journal

Proceedings of the Annual Meeting of the Cognitive Science Society, 46(0)

Authors

Shin, Minjung
Choi, Seongho
Heo, Yu-Jung
et al.

Publication Date

2024

Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

CogME: A Cognition-Inspired Multi-Dimensional Evaluation Metric for Story Understanding

Minjung Shin¹(mjshin77@snu.ac.kr), Seongho Choi²(shchoi@bi.snu.ac.kr),

Yu-Jung Heo³(yj.heo@kt.com), Minsu Lee⁴(minsulee@snu.ac.kr)

Byoung-Tak Zhang^{1,2,4}(btzhang@snu.ac.kr), Jeh-Kwang Ryu⁵(ryujk@dongguk.edu)

¹Interdisciplinary Program in Cognitive Science, Seoul National University,

²Department of Computer Science and Engineering, Seoul National University, ³KT,

⁴AI Institute of Seoul National University(AIIS), ⁵Department of Physical Education, Dongguk University

^{1,2,4} Seoul, 08826, Republic of Korea, ³ Seoul, 06763, Republic of Korea, ⁵ Seoul, 04620, Republic of Korea

Abstract

We introduce CogME, a cognition-inspired, multi-dimensional evaluation metric for AI models focusing on story understanding. CogME is a framework grounded in human thinking strategies and story elements that involve story understanding. With a specific breakdown of the questions, this approach provides a nuanced assessment revealing not only AI models' particular strengths and weaknesses but also the characteristics of the benchmark dataset. Our case study with the DramaQA dataset demonstrates a refined analysis of the model and the benchmark dataset. It is imperative that metrics align closely with human cognitive processes by comprehending the tasks' nature. This approach provides insights beyond traditional overall scores and paves the way for more sophisticated AI development targeting higher cognitive functions.

Keywords: artificial intelligence; video story understanding; video question answering; evaluation metric

Introduction

In recent years, the development of artificial intelligence (AI) models, particularly pre-trained Large Language Models (LLMs) (Vaswani et al., 2017) and diffusion models (Ho, Jain, & Abbeel, 2020), has made remarkable progress. These models have demonstrated impressive performance in creative tasks, including generating images (Rombach, Blattmann, Lorenz, Esser, & Ommer, 2022), videos (Singer et al., 2023), and narrative storytelling (Brown et al., 2020).

However, skepticism remains about their ability to 'understand,' as highlighted in recent discussions (Van Noorden & Perkel, 2023; Millièrè & Buckner, 2024). Indeed, while AI models often outperform humans in generating text and images, their performance in understanding does not reach their outstanding generative outputs (West et al., 2023). This limitation becomes prominent in multi-modal AI models, where integrating and interpreting various data forms – such as image, video, and textual information – imposes considerable challenges. In line with this limitation, video story understanding models demonstrate a significant gap compared to human story comprehension (Zhong et al., 2022).

Another key issue is the inadequacy of current evaluation metrics for AI models. There are widespread arguments that existing metrics are too general and fail to provide a comprehensive analysis of these models. These metrics often rely on aggregate scores, which can obscure true model performance and hinder understanding the benchmark dataset's detailed features used for training (Gundersen & Kjensmo, 2018; Burnell et al., 2023). This becomes increasingly pronounced in

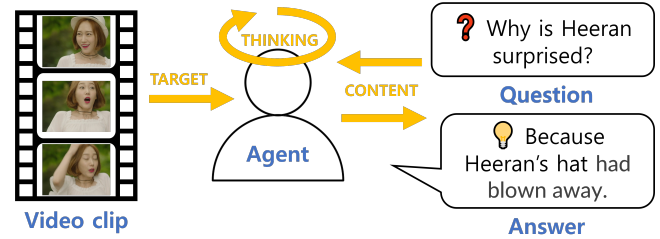


Figure 1: An illustration of CogME framework for an example of DramaQA dataset. It shows a situation in which an *Agent* predicts the *Answer* from the given *Video clip* and *Question*. Orange arrows indicate the process involves three story understanding components: *TARGET*, *CONTENT*, and *THINKING*.

more complex tasks, like understanding video stories. Consequently, there is an urgent need for a new method that aligns with human cognition for a more fine-grained assessment of AI model performance. Addressing this pressure would require a thorough understanding of the intrinsic nature of given tasks and the evaluation target.

In light of this need, we have proposed a novel evaluation metric for the story understanding model, named **Cognition-inspired Multi-dimensional Evaluation (CogME)**. CogME is designed to evaluate through a unique lens: A specific breakdown of the questions. The breakdown is grounded in multi-dimensional criteria that consider human thinking strategy and story elements.

The unique design is based on the following proposition: *If an agent answered a specific question appropriately, it means that "The agent understood the CONTENT of the TARGET through a way of THINKING." It also means "The question required the agent's THINKING about the CONTENT of the TARGET."* Our approach analyzes the context intricately to provide a richer explanation of how AI models tackle each query. This framework not only identifies the strengths and weaknesses of AI in handling various questions but also highlights the dataset's distinctive features.

We also present the results from an in-depth case study using DramaQA, a representative benchmark for video story understanding, along with its baseline model. Fig. 1 shows

how CogME is applied to DramaQA. Our findings confirm that CogME enables a thorough and systematic evaluation of both benchmark datasets and AI models.

Related Work

Narrative Comprehension of Human

Cognitive science studies on narrative comprehension, conducted through various reading comprehension tasks, have yielded seminal findings. It was found that people tend to prioritize the main aspects of a story over individual parts, indicating a higher level of cognitive engagement with narrative comprehension (Thorndyke, 1977). The concept of *schema* (Bartlett & Burt, 1933; Brewer, 1985) in memory reveals that recall is an active process influenced by personal and cultural contexts, emphasizing the critical roles of context and prior knowledge in comprehension (Bransford & Johnson, 1972; Graesser, Singer, & Trabasso, 1994). Mental representations such as *scripts* and *schemata* were denoted as crucial cognitive structures aiding comprehension and inference-making (Schank & Abelson, 1975; Rumelhart, 1980).

From these research bases, the *Situation model* and *Event-indexing model* laid the groundwork for analyzing human story comprehension, providing a structured basis for evaluation. The situation model constructs detailed mental representations encompassing events, characters, and settings (Zwaan & Radvansky, 1998). The event-indexing model emphasizes five independent dimensions of understanding when reading: time, space, character, causality, and motivation in the narrative context (Zwaan, Langston, & Graesser, 1995).

Research on understanding video stories is relatively scarce compared to reading comprehension despite the rapidly increasing importance of video narratives in our daily lives. Since reading and viewing are distinct tasks with different cognitive loads (Jajdelska et al., 2019; Cohn, 2020), it is undesirable to replicate reading comprehension structures for video comprehension (Gibson, 1979; Hochberg & Brooks, 1996).

Recent research focuses on understanding video narratives in real-world contexts. Specifically, top-down approaches involve observing brain activity, not in controlled or fragmented experimental videos, but rather in watching typical movie scenes (Baldassano, Hasson, & Norman, 2018; Song, Park, Park, & Shim, 2021). They contrast with conventional cognitive psychology studies, which primarily employ a bottom-up approach focusing on the segmented visual and audio stimuli (Tan, 2018). These integrated approaches emphasize considering both perceptual and narrative aspects to understand how people interpret video narratives in real-world situations.

Evaluation of Machine Comprehension

Machine Reading Comprehension (MRC) is the most prominent within machine comprehension (Hirschman, Light, Breck, & Burger, 1999). Question-answering (QA) has been widely adopted to evaluate text understanding of MRC

models. (Borges, 2013; Baradaran, Ghiasi, & Amirkhani, 2022). However, efforts to access the MRC models systematically have only recently been made. Dunietz et al. argued that existing MRC metrics lack clarity and could be improved by using templates derived from the definition of comprehension (Dunietz et al., 2020). To address this issue, they employed the *Event-indexing model* from human studies (Zwaan et al., 1995) to posit four elements that machines should incorporate for better reading comprehension: place, time, causality, and motivation. Similarly, Weston et al. demonstrated empirical test results of multiple AI agents' textual understanding abilities with structured QA skillsets categorized twenty types of questions for understanding and reasoning with text (Weston et al., 2015). However, these tests are limited to being constructed with fragmentary and artificial descriptions. Further research is needed before their method can be applied at the level of complexity and richness typically found in everyday human storytelling.

Video understanding models have yet to keep pace with the growing demand, even though video-based storytelling has recently emerged as one of the most prominent forms of media content. Developing AI that can understand video stories is challenging, given that it requires an all-inclusive process to analyze images, scripts, and sounds with temporal dependencies, natural language, and various levels of reasoning (Bebensee & Zhang, 2021). Despite the difficulty, several efforts to develop video understanding AI have centered on large-scale video datasets (Tapaswi et al., 2016; Lei, Yu, Bansal, & Berg, 2018; Yu et al., 2019; Garcia, Otani, Chu, & Nakashima, 2020; Choi et al., 2021; Yang, Miech, Sivic, Laptev, & Schmid, 2021), but the technologies do not extend beyond simple image processing tasks such as detecting or tracking objects.

From the evaluation perspective, the most prevalent approach involves building massive QA datasets, leveraging open-ended or multiple-choice QAs for AI training and testing (Patel, Parikh, & Shastri, 2020). However, existing evaluation methods heavily count on unidimensional metrics (Aafaq, Mian, Liu, Gilani, & Shah, 2019), such as basic QA accuracy scores for the multiple-choice QA datasets, which often fall short in providing a comprehensive explanation of the model's performance. In the case of open-ended QA, the automatic evaluation primarily depends on n-gram-based sentence similarity measures such as BLEU (Papineni, Roukos, Ward, & Zhu, 2002), METEOR (Banerjee & Lavie, 2005), ROUGE (Lin, 2004), and CIDEr (Vedantam, Zitnick, & Parikh, 2015). These metrics, focusing on the similarity between the reference and generated sentences, frequently do not align with human judgments, reflecting a significant disconnect. These automatic metrics are usually insufficient to provide a refined interpretation like the strengths and weaknesses of the AI model (Nema & Khapra, 2018).

On the other hand, there are several efforts in human evaluations to appraise the effectiveness of automatic metrics in accurately assessing AI performance. (Chen, Stanovsky, Singh,

Table 1: The sub-components within TARGET

Elements	Definition
Character	Information of individuals featured in the video.
Object	Items and body parts featured in the video.
Place	Spatial information of the story in the video.
Conversation	Characters’ dialogues, monologues, speech sounds, and text messages.
Behavior	Movements and actions of the subjects in the video.
Event	Information about what happened in the video.
Emotion	The feelings expressed by the subject in the video.
Commonsense	Concepts and knowledge that people universally accept in a given culture.

& Gardner, 2019; Garcia et al., 2020). While these efforts indicate that human evaluations are beneficial for assessing AI performance, they also reveal a significant challenge: the fine-grained evaluation is not effectively achieved in proportion to the resources and costs involved in the evaluation process.

To sum up, evaluating AI solely based on its QA accuracy or similarity to human performance is inadequate, especially for tasks with high complexity, such as video story understanding. Accordingly, a structured framework that meticulously analyzes both the nature of the understanding process and the unique features of the medium is needed. In response to this, we propose a new metric that not only reflects the existing frameworks of human story comprehension but also incorporates the distinct attributes of video storytelling.

New Evaluation Paradigm Based on the Understanding Processes of Humans

To evaluate understanding competence thoroughly, we developed multifaceted criteria integrating video narrative elements and thinking strategies involving queries. In analyzing the story elements provided by the video, we have adapted the Situation Model (Zwaan & Radvansky, 1998) and the Event-Indexing Model (Zwaan et al., 1995) and expanded them to better suit video narratives. Regarding human thinking strategies, we referred to Bloom’s Taxonomy, widely accepted as a representative framework demonstrating the hierarchical structure of cognitive processes (Bloom, 1956; Anderson & Krathwohl, 2001). Each level of Bloom’s taxonomy represents a different cognitive skill, ranging from the basic recall of facts and grasping details to reasoning hidden matter and, ultimately, evaluation and creation.

Drawing on insights from previous models in cognitive

Table 2: The sub-components within CONTENT

Elements	Definition
Identity	Personal information of subjects or names of objects in the story.
Feature	Characteristics, traits, or atmosphere of subjects and/or objects.
Relationship	The relationships between two or more targets.
Means	Instruments or methods used to achieve a particular purpose.
Context	Story-line revealed through the conversations or interactions between characters.
Sequence	Related events with time series and the changes before and after.
Causality	Causes and consequences of a particular change: natural or mechanical.
Motivation	Changes resulting from actions involving personal preferences or intentions.

Table 3: The sub-components within THINKING

Elements	Definition
Recall	Retrieving or recollecting factual information in the scene.
Grasping	Perceptions or interpretations of the scene with temporal and spatial changes.
Reasoning	Making logical judgments from circumstantial evidence not direct observations.

science, we have developed a multi-dimensional metric, CogME, which consists of three key components: TARGET, CONTENT, and THINKING. **TARGET** refers to the information perceived by watching the video, **CONTENT** to the knowledge acquired through the target information, and **THINKING** to the cognitive process of deriving knowledge from the information. The three components, representing cognitive processes, integrate to interpret the story elements presented in the video, succinctly expressed as “*I understand the CONTENT of the TARGET through a way of THINKING.*”

The sub-components of TARGET and CONTENT were determined by analyzing the story elements necessary for understanding video stories. The sub-components of THINKING were decided to align with the range required among human thinking strategies. Tables 1 – 3 demonstrate the sub-components of TARGET, CONTENT, and THINKING.

Materials and Methods

Application to VideoQA Dataset: DramaQA

This study evaluated the DramaQA dataset, which included ~16k human-generated QA pairs closely centered around the narrative of a TV series *Another Miss Oh*, along with character-level annotations (Choi et al., 2021; Bebensee &

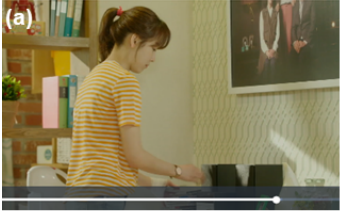

	Q: What is Haeyoung1 wearing?								
	Target	Character	Object	Place	Conversation	Behavior	Event	Emotion	Commonsense
		✓	✓						
	Content	Identity	Feature	Relationship	Means	Context	Sequence	Causality	Motivation
		✓							
	Thinking	Recall	Grasping	Reasoning					
		✓							
Q: Why does Heeran envy Haeyoung1?									
	Target	Character	Object	Place	Conversation	Behavior	Event	Emotion	Commonsense
		✓			✓			✓	
	Content	Identity	Feature	Relationship	Means	Context	Sequence	Causality	Motivation
						✓			✓
	Thinking	Recall	Grasping	Reasoning					
				✓					

Figure 2: Examples of tags applied to questions (a) Cases of tagged to questions in *shot*-level video, which require simple recall. (b) Cases of tagged to questions in *scene*-level video, which require comprehensive reasoning.

Zhang, 2021). The character-centered annotations and five-option multiple-choice QA pairs in DramaQA were generated by approximately two to five trained human annotators using a consistent manual for all 18 episodes. The dataset was designed to reflect various narrative elements in the questioning stage, from seeking simple information from the video to reasoning complex causality about the stories (Heo et al., 2019). The baseline model of DramaQA, i.e., the Multi-level Context Matching (MCM) model (Choi et al., 2021), was trained with the 1st to 12th episodes, validated with the 13th to 15th episodes and tested with the 16th to 18th episodes.

Annotating the Understanding Components

To determine what information, knowledge, and thinking strategies are required for answering questions, we annotated 4,385 questions from episodes 13th to 15th using the understanding sub-components defined in the CogME framework. Two specialists in cognitive science elaborately analyzed the given videos and questions, tagging the required sub-components in each question. To ensure consistent annotation tagging, both individuals followed the same predefined manual and resolved any discrepancies in inter-rater annotations through discussion to reach a consensus.

Fig. 2 illustrated two examples of this annotation.¹ Fig. 2(a) shows an example of a question answerable by simply recalling a single cue from a *shot*-level video. It involves identifying clothing-related information about the only person in the shot. In contrast, Fig. 2(b) illustrates a complex scenario where tagging sub-components in a *scene*-level video highlights a different level of narrative understanding from the example shown in Fig. 2(a). During a three-minute and

six-second runtime, two characters are shown walking in the park and chatting. To answer the question, the agent must not only recall and grasp the content of the conversation but also infer why Heeran said to envy Haeyoung, which is not directly mentioned in the dialog.

For THINKING strategies, we assumed that higher cognition encompasses lower ones based on the hierarchical Bloom’s taxonomy (Bloom, 1956), so we labeled only the highest thinking component. Regarding the TARGET and CONTENT of the question, we tagged up to three if multiple sub-elements were involved in a single question.

Scoring the Questions and Prediction Results

In the context of the logical complexity of the THINKING module, a weight of 2 and 3 was assigned to *grasping* and *reasoning*, respectively, assuming grasping includes the simple recall and the reasoning process retains the recall and grasping. All tagged sub-components were multiplied by the weight given to the THINKING component, as the depth of the thinking strategy determines the overall difficulty of the question. In the example shown in Fig. 2, for question in Fig. 2(a), which requires *recall*, all labeled elements are worth 1 point, while for question in Fig. 2(b), which requires *reasoning*, all labeled elements are worth 3 points.

Even identical questions can be endowed with different weights. This variance depends on the diverse contexts of the given video, ultimately affecting the complexity of the question. Considering the question, ‘Who is smiling?’ if there is only one character in the video, the prediction can be made simply by recalling who it is. However, when there are two or more characters, the prediction requires a higher level of thinking, which involves discriminating the smiling character from the others and then identifying that character. This grasping process certainly includes recalling relevant details.

The model’s accuracy for each sub-component was calculated by scoring correct predictions as 1 and incorrect ones as

¹The DramaQA dataset features memory capacity criteria related to the length of the video segments. The criteria include two types of video clip: 1) *shot*-level video without camera cut, spanning a few seconds, and 2) *scene*-level video with multiple events in a single location, spanning a few minutes

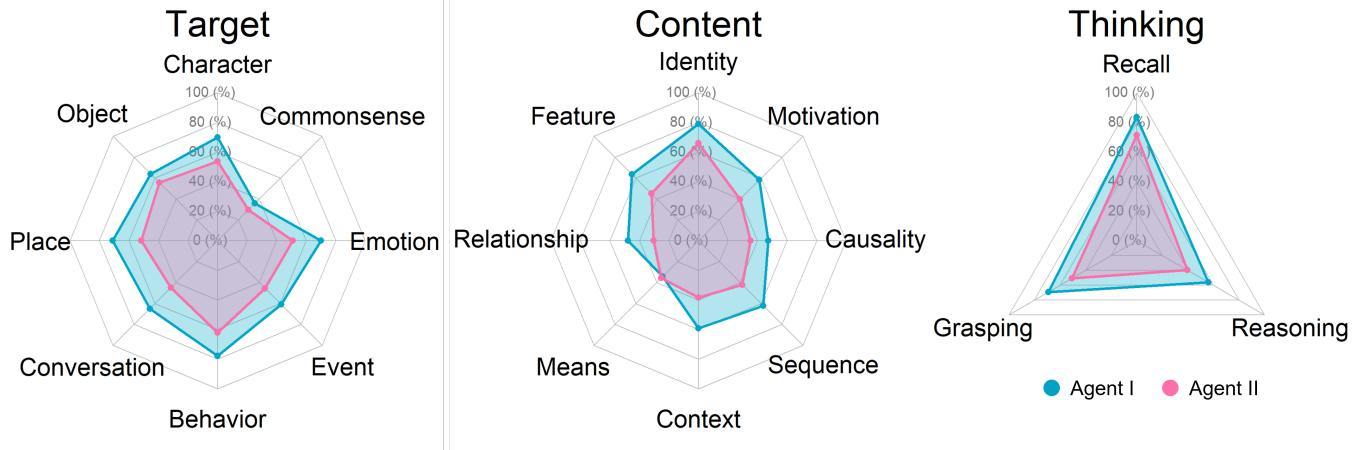


Figure 3: Performance profiles of two models. The vertex of each polygon represents the ratio (%) of correct predictions for the DramaQA dataset. Each radar plot represents TARGET(left), CONTENT(middle), THINKING(right) component. Light blue areas indicate the performance profiles of Agent I (MCM model), and pink areas display the performance profiles of Agent II (MemN2N model).

0. We then computed the overall success rate by comparing the number of correct predictions to the total attempts. To address the imbalanced frequency of each sub-component, the achievement rate was expressed as a percentage.

The analysis was conducted in R 4.1.1 (R Core Team, 2021) environment for arithmetic computing and visualization. The polygonal profiles were produced using the fmsb 0.7.1 package (Nakazawa, 2021).

Comparison of Multiple Agent Performances

In additional experiments, two agents trained with the same dataset were evaluated with CogME to compare their performances. The two models were the MCM model (Choi et al., 2021) (Agent I) and the baseline model of MemN2N (Sukhbaatar, Szlam, Weston, & Fergus, 2015) (Agent II)². The two models were examined after being trained on the same dataset, DramaQA. Their performance profiles were generated using CogME.

Results

Evaluating Model Performances

For a fine-grained analysis of model performance, the model’s accuracy for each sub-component defined by CogME was scored based on correct predictions, and each accuracy was calculated by comparing correct predictions to the total attempts.

Fig. 3 shows the multi-dimensional accuracy profile for each understanding component obtained from applying CogME to the dataset. The profile indicates that the Agents demonstrate varying levels of competence across different sub-components, as depicted by the uneven shape of the

polygons. For example, based on Agent I, Identity of CONTENT shows an accuracy of 79.1%, while Means only achieves 34.3% accuracy.

Additionally, we observed distinct disparities when comparing the results of two different models, Agent I (MCM model, represented by a light blue polygon) and Agent II (MemN2N model, represented by a pink polygon). When the two models were trained on the same dataset, the overall correct prediction rates were 73.4% for Agent I and 58.7% for Agent II, indicating a difference of 14.7%. This difference clearly illustrates that, as is often the case in many benchmark analyses, the MCM model specialized for the DramaQA dataset outperforms the MemN2N model, which primarily targets natural language processing.

However, in a detailed breakdown according to CogME’s criteria, this discrepancy is not uniform across all areas but varies by specific factors, as shown in Fig. 3. For instance, Agent II leads by approximately 1% in questions requiring the identification of the Means, while the gap extends to 20% in questions involving the Conversation. This discrepancy underscores the significance of the CogME metric in providing an in-depth understanding of each model’s performance, a potential that has yet to be fully explored in this context.

Analyzing Questions in DramaQA Dataset

Alongside performance profiling, annotations based on CogME enable us to figure out the dataset’s features in terms of data distribution. Fig. 4 shows the distribution of sub-components tagged for the questions in the DramaQA dataset, which characterizes the benchmark.

The uneven distribution of sub-components reflects an unequal consideration of aspects of story comprehension during the dataset creation phase. The prominent bars indicate that the dataset is heavily skewed towards questions that ask superficial information, like *recalling a character’s identity*.

²It should be noted that these experiments were not meant to identify the model with the better performance but only to compare them objectively. Accordingly, in the **Results** section, we refer to the models as Agent I and Agent II instead of their names

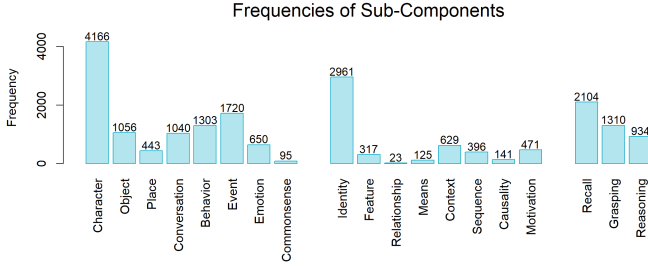


Figure 4: Frequencies of sub-components tagged in the questions of the DramaQA dataset. Each bar shows the number of times a sub-component was labeled out of 4,385 questions.

Moreover, the sub-components significantly underrepresented in the dataset tend to align with lower accuracy in model performance. Notably, all four elements that appeared less than 5% in the dataset (i.e., Commonsense, Relationship, Means and Causality) showed low accuracies that are below 50%: Commonsense (35.8%), Relationship (47.8%), Means (34.3%), Causality (47.2%).

Discussion

In this study, we introduced a novel framework, CogME, centered around the features of the posed questions. This approach is grounded in structured metrics that consider human thinking strategies and story elements using a top-down perspective. Unlike conventional AI evaluation methods, which emphasize overall scores that lead to a lack of robust evaluation, CogME provides a multi-dimensional quantified profile for AI models. This profile provides insight into the model’s strengths and weaknesses in understanding abilities by applying the metric to an existing dataset.

Our comparison of the two models using the CogME metric revealed detailed dissimilarities that their aggregate scores could not explain. These variations highlight the importance of a multi-dimensional evaluation approach for accurately assessing AI models’ capabilities. This assessment is expected to apply to both machines and humans, providing a comprehensive quantification of the agents’ levels of understanding. (Lee, Heo, Choi, Choi, & Zhang, 2023).

Furthermore, CogME’s fine-grained evaluation not only assesses the AI models but also offers an analysis of the benchmark dataset, providing deeper insights into the models’ capabilities. For instance, the observed link between low frequency and low accuracy³, indicates that learning deficiencies can impact QA performance, highlighting the need for a more balanced dataset covering various aspects of narrative comprehension. Moreover, as noted in the **Narrative Comprehension of Human** section, people focus on a story’s central aspects rather than individual instances (Thorndyke,

³Although we only analyzed questions from the validation set, we assumed that the CogME profile of the training set would be similar based on the preliminary analysis that demonstrated a similar distribution of question types across the datasets. (Choi et al., 2021)

1977). However, our analysis reveals that this dataset predominantly collected fragmentary information instead of emphasizing the story’s central or structural elements.

This insight leads to the establishment of a CogME framework as a guideline for designing new datasets. Generating a massive QA dataset makes it challenging to ensure sufficient variety in question types and sub-components, as seen in Fig. 4. These maldistribution issues have been noted not only in DramaQA but also in many other QA datasets (Garcia-Molina, Joglekar, Marcus, Parameswaran, & Verroios, 2016). In this context, the CogME framework could serve as a theoretical foundation for proper data allocation in datasets, whether through crowdsourcing or automatic question generation.

We acknowledge that a challenge in our study is that sub-components were tagged manually in the provided questions and videos. It was inevitable to capture the elements comprehensively, as even the identical questions can vary in different video contexts (see **Scoring the Questions and Prediction Results** section). Despite being cumbersome, manual annotation ensures accurate evaluation by aligning with the nuanced content of the videos and related queries. However, in the future, using a multi-modal classification model to automatically annotate sub-components in CogME could streamline the evaluation process. Such automation would not only simplify the evaluation process but also allow for the scalability of larger and more complex datasets.

Additionally, scoring multiple-choice questions can result in some information loss. According to our annotating, even if the model correctly recognized the Character’s Identity but failed to infer other information, like Emotion, it would score zero for that question, including the Character’s. By incorporating rubric-like methods used in pedagogy (Brookhart, 2018), we argue that this metric could be adapted to other tasks like open-ended or fill-in-the-blank tests, summaries, and rewriting, which can be analyzed through understanding sub-components (Lee et al., 2023).

In conclusion, this study introduces the CogME framework, offering a multi-dimensional analysis focusing on story understanding that surpasses the limitations of overall scores, like accuracy rates. CogME’s potential extends to various AI tasks and dataset designs, suggesting its adaptability and utility in advancing AI assessment toward more nuanced and sophisticated dimensions. This work also establishes a new benchmark, paving the way for more comprehensive approaches to developing AI agents.

Acknowledgments

We deeply thank the reviewers for providing kind and helpful comments. This work was partly supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grants (No. 2017-0-01772, VTT — No. 20220-00951, LBA) and by the National Research Foundation of Korea (NRF) grant (No. RS-2024-00358416, AutoRL) funded by the Korea Government (MSIT).

References

- Aafaq, N., Mian, A., Liu, W., Gilani, S. Z., & Shah, M. (2019). Video description: A survey of methods, datasets, and evaluation metrics. *ACM Computing Surveys (CSUR)*, 52(6), 1–37. doi: 10.1145/3355390
- Anderson, L. W., & Krathwohl, D. R. (Eds.). (2001). *A taxonomy for learning, teaching, and assessing: A revision of bloom's taxonomy of educational objectives*. New York: Longman.
- Baldassano, C., Hasson, U., & Norman, K. A. (2018). Representation of real-world event schemas during narrative perception. *Journal of Neuroscience*, 38(45), 9689–9699.
- Banerjee, S., & Lavie, A. (2005). Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization* (pp. 65–72).
- Baradaran, R., Ghiasi, R., & Amirkhani, H. (2022). A survey on machine reading comprehension systems. *Natural Language Engineering*, 28(6), 683–732.
- Bartlett, F., & Burt, C. (1933). Remembering: A study in experimental and social psychology. *British Journal of Educational Psychology*, 3(2), 187–192.
- Bebensee, B., & Zhang, B.-T. (2021). Co-attentional transformers for story-based video understanding. In *Icassp 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (p. 4005–4009). doi: 10.1109/ICASSP39728.2021.9413868
- Bloom, B. S. (1956). Taxonomy of educational objectives: The classification of educational goals. In B. S. Bloom, M. D. Engelhart, E. Furst, W. H. Hill, & D. R. Krathwohl (Eds.), *Handbook I: Cognitive domain*. David McKay Company, Inc.
- Bransford, J. D., & Johnson, M. K. (1972). Contextual prerequisites for understanding: Some investigations of comprehension and recall. *Journal of verbal learning and verbal behavior*, 11(6), 717–726.
- Brewer, W. F. (1985). The story schema: Universal and culture-specific properties. In D. Olson, N. Torrance, & A. Hildyard (Eds.), *Literacy, language and learning: the nature and consequences of reading and writing* (pp. 167–194). Cambridge: Cambridge University Press Cambridge, UK.
- Brookhart, S. M. (2018). Appropriate criteria: Key to effective rubrics. *Frontiers in Education*, 3. doi: 10.3389/educ.2018.00022
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... Amodei, D. (2020). Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, & H. Lin (Eds.), *Advances in neural information processing systems* (Vol. 33, pp. 1877–1901). Curran Associates, Inc.
- Burges, C. J. (2013). *Towards the machine comprehension of text: An essay* (Tech. Rep. No. MSR-TR-2013-125). Redmond, WA 98052: Microsoft Research.
- Burnell, R., Schellaert, W., Burden, J., Ullman, T. D., Martinez-Plumed, F., Tenenbaum, J. B., ... others (2023). Rethink reporting of evaluation results in ai. *Science*, 380(6641), 136–138. doi: 10.1126/science.adf6369
- Chen, A., Stanovsky, G., Singh, S., & Gardner, M. (2019). Evaluating question answering evaluation. In *Proceedings of the 2nd workshop on machine reading for question answering* (pp. 119–124). doi: 10.18653/v1/D19-5817
- Choi, S., On, K.-W., Heo, Y.-J., Seo, A., Jang, Y., Lee, M., & Zhang, B.-T. (2021, May). Dramaqa: Character-centered video story understanding with hierarchical qa. In *Proceedings of the aaai conference on artificial intelligence* (Vol. 35, p. 1166–1174). doi: 10.1609/aaai.v35i2.16203
- Cohn, N. (2020). Your brain on comics: A cognitive model of visual narrative comprehension. *Topics in cognitive science*, 12(1), 352–386.
- Dunietz, J., Burnham, G., Bharadwaj, A., Rambow, O., Chu-Carroll, J., & Ferrucci, D. (2020, July). To test machine comprehension, start by defining comprehension. In *Proceedings of the 58th annual meeting of the association for computational linguistics* (pp. 7839–7859). doi: 10.18653/v1/2020.acl-main.701
- Garcia, N., Otani, M., Chu, C., & Nakashima, Y. (2020, Apr.). Knowit vqa: Answering knowledge-based questions about videos. *Proceedings of the AAAI Conference on Artificial Intelligence*, 10826–10834. doi: 10.1609/aaai.v34i07.6713
- Garcia-Molina, H., Joglekar, M., Marcus, A., Parameswaran, A., & Verroios, V. (2016). Challenges in data crowdsourcing. *IEEE Transactions on Knowledge and Data Engineering*, 28(4), 901–911. doi: 10.1109/TKDE.2016.2518669
- Gibson, J. J. (1979). *The ecological approach to visual perception*. Houghton, Mifflin and Company. doi: 10.4324/9781315740218
- Graesser, A. C., Singer, M., & Trabasso, T. (1994). Constructing inferences during narrative text comprehension. *Psychological Review*, 101(3), 371–95. doi: 10.1037/0033-295x.101.3.371
- Gundersen, O. E., & Kjensmo, S. (2018, Apr.). State of the art: Reproducibility in artificial intelligence. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1). doi: 10.1609/aaai.v32i1.11503
- Heo, Y.-J., On, K.-W., Choi, S., Lim, J., Kim, J., Ryu, J.-K., ... Zhang, B.-T. (2019). Constructing hierarchical q&a datasets for video story understanding. In *Aaai 2019 spring symposium* (pp. 1–9).
- Hirschman, L., Light, M., Breck, E., & Burger, J. D. (1999). Deep read: A reading comprehension system. In *Proceedings of the 37th annual meeting of the association for computational linguistics* (pp. 325–332). doi: 10.3115/1034678.1034731
- Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, & H. Lin (Eds.), *Advances in neural information processing systems* (Vol. 33, pp. 6840–6851).

- Curran Associates, Inc.
- Hochberg, J., & Brooks, V. (1996). The perception of motion pictures. *Cognitive Ecology*, 205–292. doi: 10.1016/B978-012161966-4/50008-6
- Jajdelska, E., Anderson, M., Butler, C., Fabb, N., Finnigan, E., Garwood, I., ... others (2019). Picture this: a review of research relating to narrative processing by moving image versus language. *Frontiers in Psychology*, 10, 1161.
- Lee, M., Heo, Y.-J., Choi, S., Choi, W. S., & Zhang, B.-T. (2023). Video turing test: A first step towards human-level ai. *AI Magazine*, 44(4), 537–554. doi: 10.1002/aaai.12128
- Lei, J., Yu, L., Bansal, M., & Berg, T. L. (2018). Tvqa: Localized, compositional video question answering. In *Empirical methods in natural language processing*.
- Lin, C.-Y. (2004). Rouge: A package for automatic evaluation of summaries. In *Proc. of workshop on text summarization branches out, post conference workshop of acl 2004* (pp. 74–81).
- Millière, R., & Buckner, C. (2024). A philosophical introduction to language models—part i: Continuity with classic debates. *arXiv preprint arXiv:2401.03910*.
- Nakazawa, M. (2021). fmsb: Functions for medical statistics book with some demographic data [Computer software manual]. Retrieved from <https://CRAN.R-project.org/package=fmsb> (R package version 0.7.1)
- Nema, P., & Khapra, M. M. (2018). Towards a better metric for evaluating question generation systems. In *Proceedings of the 2018 conference on empirical methods in natural language processing* (pp. 3950–3959).
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the association for computational linguistics* (pp. 311–318). doi: 10.3115/1073083.1073135
- Patel, D., Parikh, R., & Shastri, Y. (2020). Recent advances in video question answering: A review of datasets and methods. In A. D. Bimbo et al. (Eds.), *Pattern recognition. icpr international workshops and challenges - virtual event, january 10-15, 2021, proceedings, part ii* (Vol. 12662, p. 339–356). Springer. doi: 10.1007/978-3-030-68790-8_27
- R Core Team. (2021). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 10684–10695).
- Rumelhart, D. E. (1980). Schemata: The building blocks of cognition. In B. W. Spiro RJ Bruce BC (Ed.), *Theoretical issues in reading comprehension* (pp. 33–58). Hillsdale, NJ: Lawrence Erlbaum Associates, Inc.
- Schank, R. C., & Abelson, R. P. (1975). Scripts, plans, and knowledge. In *Proceedings of the 4th international joint conference on artificial intelligence - volume 1* (p. 151–157). Morgan Kaufmann Publishers Inc.
- Singer, U., Polyak, A., Hayes, T., Yin, X., An, J., Zhang, S., ... Taigman, Y. (2023). Make-a-video: Text-to-video generation without text-video data. In *The eleventh international conference on learning representations*. Retrieved from <https://openreview.net/forum?id=nJfylDvgz1q>
- Song, H., Park, B.-y., Park, H., & Shim, W. M. (2021). Cognitive and neural state dynamics of narrative comprehension. *Journal of Neuroscience*, 41(43), 8972–8990. doi: 10.1523/jneurosci.0037-21.2021
- Sukhbaatar, S., Szlam, A., Weston, J., & Fergus, R. (2015). End-to-end memory networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 28). Curran Associates, Inc., NY, USA.
- Tan, E. S. (2018). A psychology of the film. *Palgrave Communications*, 4(1), 1–20. doi: 10.1057/s41599-018-0111-y
- Tapaswi, M., Zhu, Y., Stiefelhagen, R., Torralba, A., Urtasun, R., & Fidler, S. (2016, June). Movieqa: Understanding stories in movies through question-answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition (cvpr)* (p. 4631–4640).
- Thorndyke, P. W. (1977). Cognitive structures in comprehension and memory of narrative discourse. *Cognitive Psychology*, 9(1), 77–110. doi: 10.1016/0010-0285(77)90005-6
- Van Noordén, R., & Perkel, J. M. (2023). Ai and science: what 1,600 researchers think. *Nature*, 621(7980), 672–675.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). Attention is all you need. In I. Guyon et al. (Eds.), *Advances in neural information processing systems* (Vol. 30). Curran Associates, Inc.
- Vedantam, R., Zitnick, C. L., & Parikh, D. (2015). Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (p. 4566–4575). IEEE Computer Society.
- West, P., Lu, X., Dziri, N., Brahman, F., Li, L., Hwang, J. D., ... others (2023). The generative ai paradox: “what it can create, it may not understand”. In *The twelfth international conference on learning representations*.
- Weston, J., Bordes, A., Chopra, S., Rush, A. M., van Merriënboer, B., Joulin, A., & Mikolov, T. (2015). Towards ai-complete question answering: A set of prerequisite toy tasks. *arXiv preprint arXiv:1502.05698*.
- Yang, A., Miech, A., Sivic, J., Laptev, I., & Schmid, C. (2021, October). Just ask: Learning to answer questions from millions of narrated videos. In *Proceedings of the IEEE/CVF international conference on computer vision (iccv)* (p. 1686–

- 1697).
- Yu, Z., Xu, D., Yu, J., Yu, T., Zhao, Z., Zhuang, Y., & Tao, D. (2019, Jul.). Activitynet-qa: A dataset for understanding complex web videos via question answering. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01), 9127-9134. doi: 10.1609/aaai.v33i01.33019127
- Zhong, Y., Ji, W., Xiao, J., Li, Y., Deng, W., & Chua, T.-S. (2022, December). Video question answering: Datasets, algorithms and challenges. In Y. Goldberg, Z. Kozareva, & Y. Zhang (Eds.), *Proceedings of the 2022 conference on empirical methods in natural language processing* (pp. 6439–6455). Abu Dhabi, United Arab Emirates: Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.432
- Zwaan, R. A., Langston, M. C., & Graesser, A. C. (1995). The construction of situation models in narrative comprehension: An event-indexing model. *Psychological Science*, 6(5), 292-297. doi: 10.1111/j.1467-9280.1995.tb00513.x
- Zwaan, R. A., & Radvansky, G. A. (1998). Situation models in language comprehension and memory. *Psychological bulletin*, 123(2), 162–85.