
Supplementary Material for CLEVRER-Humans: Describing Physical and Causal Events the Human Way

Jiayuan Mao*
MIT

Xuelin Yang*
Stanford University

Xikun Zhang
Stanford University

Noah D. Goodman
Stanford University

Jiajun Wu
Stanford University

1 We bear all responsibility in case of violation of rights. The data created or used during this study are
2 openly available on the project’s website (<https://sites.google.com/stanford.edu/cleverer-humans>). We
3 confirm that the data is under CC0 license. We will provide maintenance to the website and dataset
4 regularly and upon request.

5 The rest of this supplementary document is organized as the following. First, in Section **A** we
6 provide visualizations of data collected in CLEVRER-Humans, as well as more analysis on the
7 comparison between human causal judgments and heuristics-based labels. In Section **B**, we describe
8 the implementation details of models studied in the main paper and add additional failure case
9 analysis of models. Next, in Section **C**, we describe the user interface for dataset collection. Finally,
10 in Section **D**, we supplement dataset sheets for CLEVRER-Humans.

11 **A Dataset Visualization and Analysis**

12 Fig. 2 and Fig. 3 show the example graph collected in the stage I (causal event cloze) and stage II
13 (binary CEG labelling), respectively. First, Fig. 2 shows that the causal cloze tasks can progressively
14 collect a large number of human-written event descriptions by re-using the response of previous
15 annotators. On average, we can obtain 29.4 descriptions per video, highlighting the advantage of our
16 design. Second, the condensed CEGs contains high-quality causal relations of physical events, as
17 shown in Fig. 3. It demonstrates both the language diversity and the richness of causal relations in
18 the CEGs of CLEVRER-Humans. These figures provide a straightforward illustration of our data
19 collection pipeline and the quality of our data.

20 **A.1 Dataset Statistics**

21 First, CLEVRER-Humans contains dense annotations of causal relations between physical events.
22 Fig. 1a and Fig. 1b show the distributions of the number of nodes and edges in each CEG. The
23 average number of CEG nodes is 4.71 and the average number of labeled edges is 12.7. These dense
24 annotations of CEGs form the rich and complicated causal structures in our dataset.

25 Second, CLEVRER-Humans offers diverse free-form language descriptions while retaining balances
26 in object properties. Fig. 1c shows the length distribution of event descriptions: the average length is
27 7.00 (as a reference, the average event description length of CLEVRER is 8.93). CLEVRER-Humans
28 has a vocabulary length of size 219, which is much greater than CLEVRER (82). Fig. 1d and Fig. 1e

*indicates equal contribution. Correspondence to: jiayuanm@mit.edu.

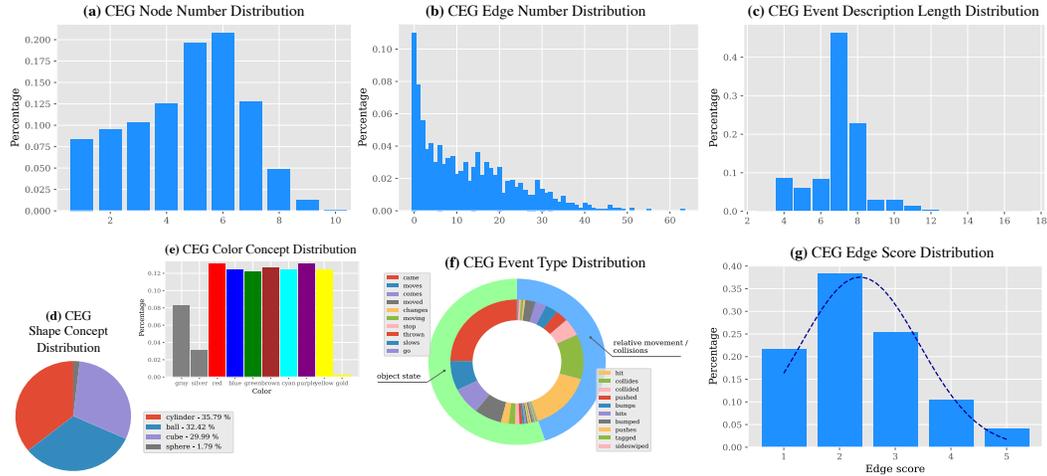


Figure 1: Statistics on the CLEVRER-Humans dataset. From left to right, the first row figures are distributions of (a) the number of nodes per CEG, (b) the number of edges per CEG, and (c) sentence lengths excluding the "which of the following is responsible for" prefix. The second row figures are distributions of (d) object shapes, (e) colors, (f) event type attributions based on verbs, and (g) CEG edge scores labelled by MTurkers, respectively.

29 show the distribution of object property concepts: colors and shapes. They remain unbiased when
 30 considering the synonyms such as "ball" and "sphere" and "gray" and "silver."

31 Next, most importantly, CLEVRER-Humans engage a variety of physical events for causal reasoning
 32 tasks. In particular, Fig. 1f shows the distribution of event types computed based on the main verb of
 33 the event description. The outer circle represents the general event families. The corresponding inner
 34 breakdowns display more than 10 variations of the expression based on verbs for each event type.
 35 In comparison, the original CLEVRER dataset contains only three event types (and verbs): enter,
 36 exit, and collide. Therefore, CLEVRER-Humans significantly improves the diversity and brings in a
 37 challenge for machines to recognize and ground these events in practice.

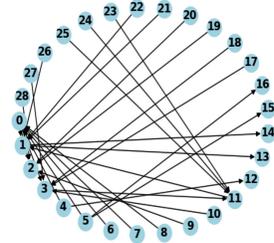
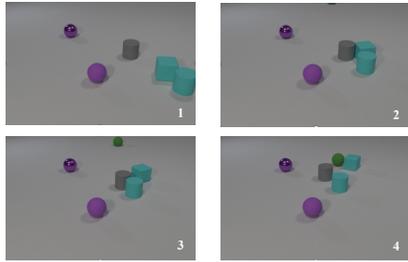
38 In the following box, we list all verbs that have been annotated by human annotators and generated
 39 by our machine generative model. We have lemmatized all verbs to remove the tense.

come, move, change, stop, throw, slow, go, travel, begin, spin, roll, stand, halt, roll, lose,
 leave, head, want, hurl, enter, hit, collide, push, bump, push, tag, sideswipe, bounce, strike,
 touch, cause

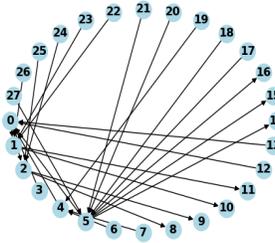
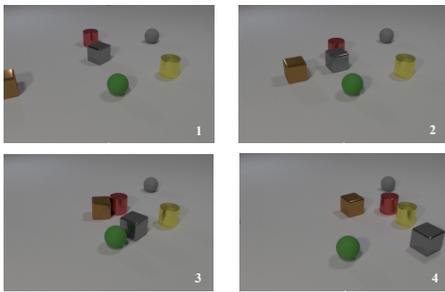
40 We also would like to point out that for some verbs, if they seem to be synonyms (e.g., bump and
 41 sideswipe), they can have subtle differences in physical grounding. For example, A bumps into B
 42 usually implies that A is moving faster than B and its collision changed the state of B. Furthermore,
 43 different tense of the same verb have different meanings in sentences: "the event that ball A moved is
 44 responsible for the collision" is different from "the event that ball A is moving is responsible for the
 45 collision." In the former case, ball A does not have to be moving while the collision happens.

46 It is possible to hand-craft a lot of rules to handle each individual cases (e.g., bump, sideswipe, roll),
 47 but that will require additional hyperparameters for thresholds, and may be hard to align with human
 48 perception.

49 Finally, CLEVRER-Humans' annotation reflects the subjective judgment of causality in physical
 50 events. CLEVRER-Humans offers 5 choices when asking MTurkers to label the causality level.
 51 Fig. 1g shows the distribution of edge scores with an average of 2.37. Note that this distribution is
 52 skewed towards lower scores (as shown by the Gaussian approximation in the dotted curve). This
 53 reflects the fact that most event pairs do not have causal relationships. Finally, although we have



- 0: the cyan cube sideswipes the grey cylinder
- 1: the green ball hits the cyan cube
- 2: the green ball touches the grey cylinder
- 3: the green ball bumps into the cyan cube
- 4: the cyan cube was moving in direction opposite to the grey cylinder
- 5: the cyan cylinder came from below
- 6: the cyan cylinder struck the cyan cube
- 7: the cyan cube was struck by the cyan cylinder
- 8: the cyan cylinder moved
- 9: the green ball was pushed by the cyan cube
- 10: the gray cylinder moved
- 11: the green ball came from above
- 12: the square cyan was hit by the cyan cube causing the grey cylinder to move away
- 13: the cyan cube is pushed to the right side
- 14: the cyan cube spins counterclockwise
- 15: the green ball collided with the cyan cube
- 16: grey cylinder moved forward
- 17: the blue rubber cylinder bumped the cyan cube in the way of the green ball
- 18: the green ball slides down into the scene into the cyan cube
- 19: the green ball touched the cyan cube
- 20: the green ball bounces off the cyan cube and touches the grey cylinder
- 21: cyan is on a direct path to green ball path
- 22: the cyan cylinder hits the cyan cube
- 23: the green ball was already moving downwards
- 24: the green ball had velocity
- 25: the green ball rolled down into the screen and bounced into the cyan cube
- 26: the cyan cylinder hit the cyan cube
- 27: the green ball and the cyan cube were heading to the same point
- 28: the green ball bounces off the blue square which was traveling in the opposite direction



- 0: the brown cube starts to spin clockwise
- 1: the brown cube hits the silver cube and then hit the red cube
- 2: the brown cube came from the west
- 3: the brown cube starts to spin clockwise because gold metal sphere to swipe
- 4: the brown cube hits the metal purple cube
- 5: the brown cube hit the ash cube by the side
- 6: the metal purple cube was in the way
- 7: the red cylinder collides with the metal purple cube
- 8: the grey cube touched the green ball
- 9: the grey cube moved south
- 10: the red cube ends up in the east
- 11: the silver cube hits the green ball
- 12: the brown cube collided with the silver cube after the silver cube was struck by the red cylinder
- 13: the brown cube collided with the grey cube
- 14: the ash cube moves to the right
- 15: the gray cube spun
- 16: the silver cube hit the green ball
- 17: the red cylinder collided with the yellow cylinder
- 18: the red cylinder pushed the ash cube away from the brown cube's path
- 19: the red cylinder hit the metal purple cube
- 20: the brown cube came from the left
- 21: the red cylinder hit the ash cube from the top
- 22: ricochets off to left and bounces into red cube
- 23: the brown cube came from the side
- 24: the brown cube started off on the west and was moving east
- 25: the brown cube was pushed by an unknown and unseen force
- 26: the silver cube absorbed the force from the brown cube in order to hit the red cube
- 27: the red cylinder hit the grey cylinder which changed its direction

Figure 2: Visualization of two samples of annotations collected in Stage I causal cloze tasks. They are collected progressively by feeding the response of a user as the input of another one. The black arrows indicate the annotation orders.

54 binarized the edge labels for the sake of consistency with CLEVRER, the raw score-based judgment
 55 can be potentially helpful in other tasks such as cognitive science studies.
 56 Therefore, we can conclude that CLEVRER-Humans is a high-quality causal relation dataset with
 57 significantly more diverse event types and language descriptions than CLEVRER.

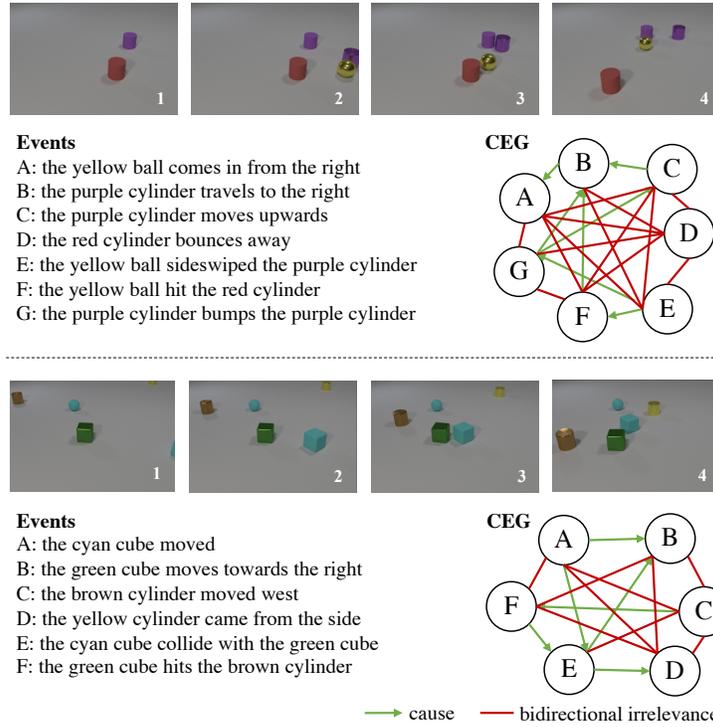


Figure 3: Visualization of two samples of CEGs in CLEVRER-Humans. The green arrows represent causal relations and the red edges represent bidirectional irrelevance. We can see the rich causal relations among physical events presented in the CEGs.

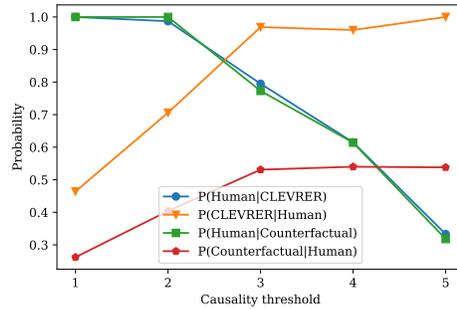


Figure 4: Effect of different causality thresholds on the binarized human causal relation. The x-axis is the ablation threshold (i.e., 4 means a score ≥ 4 represents a causal relation). The y-axis is the conditional probability.

	$Y = y_1$	$Y = y_2$	$Y = y_1 \wedge y_2$	$Y = y_1 \vee y_2$	$Y = y_1 \oplus y_2$
$P(X = \text{Human} Y)$	0.62	0.61	0.23	0.34	0.34
$P(Y X = \text{Human})$	0.96	0.54	0.29	0.62	0.33

Table 1: Comparison between different combinations of heuristics-generated causal labels and human labels, on a sampled subset of CLEVRER [1]. The entry $P(X|Y)$ denotes the fraction of event relations that are annotated as causal by protocol X given that the relations are annotated as causal by protocol Y. y_1, y_2 denote the existence of causal relations defined CLEVRER’s heuristic and Counterfactual causal relation, respectively.

58 **A.2 Comparison between Heuristic and Human Causal Judgments**

59 We supplement the effect of different thresholds on the graded causal relation in Fig. 4. In the human
60 performance study, we asked the participants to choose a threshold from 1-5 if they had to binarize
61 their judgment. The average threshold suggested by the participants is 3.6. In practice, we choose a
62 threshold of 4 to obtain the causal relation that humans are more certain about.

63 Having shown the two common heuristics-generated causal labels (CLEVRER’s and counterfactual
64 intervention) diverges from human judgment, we also provide the results on comparisons between
65 different combinations of heuristics-generated causal labels and human judgments. We use the logic
66 operators and (\wedge), or (\vee), xor (\oplus). As shown in Table 1, none of these combinations can give an
67 close enough approximation to human judgment, which further justifies our motivation to use human
68 labeled causal data for CLEVRER-Humans.

69 B Implementation Details

70 In this section, we present the implementation details of our neural-network-based event description
71 generator, the baseline models studied in the main paper, and the error bars for models across different
72 random seeds.

73 B.1 Stage II Implementation

74 We first describe the input pre-processing for neural event generators. For each object, we concatenate
75 the one-hot encoding of physical properties (including shapes, colors, and materials) and the motion
76 information (including location, orientation, velocity, angular velocity, and whether the object is
77 inside the camera view) in each of the 128 frames in a video. For each object, at each time step, the
78 input dimension is 24.

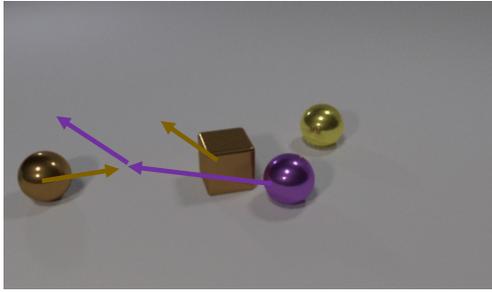
79 Our rule-based event detector for object pairs works as the following. For object pairs, we first extract
80 all segments that are composed of consecutive frames when two objects are close to each other.
81 Specifically, we say two objects are close if the L_∞ norm of the displacement vector between two
82 objects is smaller than 0.5 meter (i.e., their x, y, z displacements are all smaller than 0.5 meter).
83 Within each segment, the event detector predicts event types including moving together, object
84 approaching, and collision, based on changes in the motion information. For example, if two objects
85 are physically close for more than 20 frames without rapid changes in velocity, we consider them
86 relatively static, thus “moving together.” If both objects change directions within their close period,
87 we consider a collision happened. We can further distinguish the changes in relative positions (either
88 “bouncing back” or “one approaching another”) by the sign of the dot product of velocity vectors. For
89 any object pair, if no events are detected in the course of the entire video, we do not include this pair
90 for future captioning.

91 After getting the input sequences, we use neural event generators consisting of an encoder and a
92 decoder to produce captions. The encoder uses a linear layer and a GRU unit to encode the input
93 sequence [2]. The decoder applies Softmax on the embedding of input and the hidden state to produce
94 the attention weights. It then uses GRU and a linear layer to produce an English caption of specific
95 objects in the video. Single-object and pairwise captioning models share the same architecture but
96 are trained independently. The hidden dimension of both the encoder and the decoder is 256 for
97 single-object models and 128 for pairwise models. The dropout rate for the decoder is set to 0.3 for
98 single-object models and 0.1 for pairwise models. All models are trained with a learning rate of 0.001.
99 For the grammar check module in the post processor, we drop the sentences with two consecutively
100 repeated words. We also exclude the sentences that miss verbs or verb arguments, such as sentences
101 ending with words “from,” “to,” “at,” “is,” etc.

102 B.2 Baseline Implementation

103 **Language-Only models.** For the language-only models, we use a LSTM [3] with GloVe [4] word
104 embedding. The hidden dimension is 512 and the dropout rate is 0.2. We use the Adam optimizer [5]
105 with a learning rate of 4×10^{-4} and a weight decay of 10^{-5} . The batch size is 4. Following the data
106 splits in CLEVRER, we split 20% of the training pairs as the validation set and choose the model
107 with the highest validation accuracy.

108 **CNN+LSTM.** For the CNN+LSTM models, we use a pre-trained ResNet-50 to extract 2,048-
109 dimensional features from the video frames [6]. We uniformly sample 25 frames for each video as
110 input. The word embedding for questions is initialized by the GloVe [4] word embedding. Both
111 LSTMs (the question encoder and the video sequence encoder) have 1 layer with a hidden dimension
112 of 512. We apply a dropout rate of 0.2 on the input layer and 0.5 on the hidden layers. We use
113 the Adam [5] optimizer with a weight decay of 5×10^{-4} . The learning rate is 10^{-5} training from
114 scratch and 10^{-3} for finetuning. The batch size is 128 for both trained-from-scratch and finetuning
115 experiments. We split 20% of the training pairs as the validation set and choose the model with the
116 highest validation accuracy.

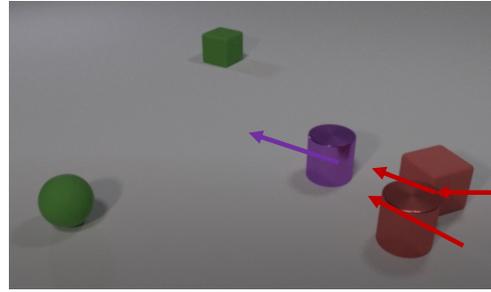


Question: What is responsible the purple ball collides with the brown cube?

Choice: The purple ball comes up.

Answer: Wrong

Model prediction: Correct



Question: What is responsible for the red cube sideswiped the purple cylinder?

Choice: The red cube bumped the red cylinder.

Answer: Wrong

Model prediction: Correct

Figure 5: Examples of common prediction errors. The arrows in the image represent the moving direction of objects of interest in the video. **Left:** failure caused by nuances in human language. While the purple ball is constantly moving upwards coordinate-wise, humans understand the phrase "comes up" as more of the later part of the trajectory (after the purple ball collides with the brown ball). Therefore, machines cannot give a correct prediction. **Right:** failure in bridging the domain shift. Humans may consider the change in the trajectory to be minor and appears not to be a deciding factor of the outcome event, but the model predicts it as a cause following similar heuristics in CLEVRER.

117 **BERT+LSTM** We supplement CNN+BERT models as a model with a stronger text encoder. The
 118 CNN is the same as in the CNN+LSTM baseline. We use the pretrained BERT uncased base model
 119 from HuggingFace library [7]. The BERT tokenizer is set to max length 32 padding and truncation.
 120 During training, We fix weights of the text encoder. We use the Adam optimizer with a weight decay
 121 of 5×10^{-4} , and a learning rate of 10^{-5} training from scratch and 10^{-3} for finetuning. The batch
 122 size is 128. We choose the model with highest validation accuracy with 20% of the training set as the
 123 validation set.

124 **ALOE.** We implement our model based on the publicly released code [8]. Since the public release
 125 does not contain training code, we implement the training procedure using the following settings.
 126 For object embeddings, we use the pre-trained MONet embeddings released by the authors. For
 127 optimization, we use the Adam [5] optimizer with a weight decay of 10^{-3} (we have also benchmarked
 128 10^{-2} , 10^{-3} and 10^{-4}). We split 5% of the training pairs as the validation set and choose the model
 129 with the highest validation accuracy.

130 B.3 Error analysis

131 We summarize the common failures of the models: for pretrained-only models, the common error
 132 comes from the failure of incorporating more diverse language and events. For example, as shown in
 133 Table 2, the program parser of NS-DR and VRDP fails to generate proper programs for descriptions
 134 in CLEVRER-Humans. The deficiency of language understanding often leads to wrong predictions.

135 For training from scratch models, one possible reason to the test errors is the nuances in human
 136 language. Specifically, models do not only need to identify the objects being referred to but also
 137 their physical properties: the cause of "the red cube slows down" can be hard to identify because
 138 speed does not appear to be as explicit as other properties such as colors and shapes. As shown
 139 in our comparison between human judgement and heuristics-based causal judgements, the nuance
 140 in language can influence human judgments, posing difficulties for machines to ground the events
 141 and simulate the reasoning process. For instance, the left figure in Fig. 5 illustrates the nuances
 142 in language resulting a discrepancy between human judgment and prediction. Moreover, for large

Event	Parsed program
The purple sphere slows down from the right.	["events", "objects", "purple", "filter_color", "sphere", "filter_shape", "unique", "filter_collision", "objects", "unique", "filter_color", "sphere", "filter_shape", "unique", "filter_collision", "unique"]
The red ball comes to a stop.	["events", "objects", "red", "filter_color", "unique", "filter_collision", "objects", "red", "filter_color", "unique", "filter_collision", "unique"]
The yellow cube comes from the right side at a fast speed.	["events", "objects", "yellow", "filter_color", "cube", "filter_shape", "unique", "filter_out", "unique"]

Table 2: Examples of errors produced by program parser. In the first row, the model cannot identify the event "slow down from the right" and gives incorrect parsing to find another object involved in a collision ("filter_collision"). In the second row, the model cannot represent the event "come to a stop" due to the expand in vocabulary and gives an incorrect output ("filter_collision"). In the third row, the model mistakenly represents the enter event as the exit event ("filter_out") because the description is more complicated in CLEVRER-Humans. We follow the notation of programs as in NSDR and VRDP.

143 models such as ALOE, learning to simulate human reasoning process from scratch based on very
 144 little data can be difficult, especially with limited training size.

145 For finetuned models, we have not seen significant improvement brought by the pretraining phase.
 146 This is primarily because of the domain gap between human judgement and heuristics-based labelling.
 147 Specifically, our human experiments have shown that $p(\text{Human} \mid \text{CLEVRER-Heuristic}) = 0.62$. That
 148 is, only 62% of the event pairs that have been labelled as causal in CLEVRER, are labelled as causal
 149 by human annotators. The right figure in Fig. 5 gives an example of the error caused by the domain
 150 shift. Future work may consider other ways of pretraining, such as pretraining on event recognition,
 151 which may be more transferable, and pretraining with other types of heuristics.

152 **C Labeling Interface**

153 We develop labeling interfaces based on boto3 with Amazon MTurk python API. We include example
154 trials of both the causal cloze tasks and CEG annotation tasks in Fig. 6 and Fig. 7, respectively.
155 The full instruction texts are provided on the labeling page of our project’s website. The estimated
156 hourly pay to the Mechanical Turk participants is about \$6.1 and the total amount spent on participant
157 compensation is about \$3500. Specifically, the cloze tests and part of the pairwise causal relationship
158 annotations were completed by users from the U.S., and the pay was \$7.7/hour (above federal
159 minimum wage). At a later stage of our project, we were unfortunately constrained by the budget
160 available to us and opened the tasks to workers outside the U.S. Thus, overall, our average hourly
161 wage is \$6.1. Our goal has always been to commit to best practice and offer fair pay to users whenever
162 possible, and we will continue to do so in the future.

163 In causal cloze tasks, the participants are asked to write an event description given a cause or outcome
164 event, as shown in Fig. 6. We specify the expectation of responses (such as using complete sentences,
165 avoiding ambiguous third-person pronouns, etc.) in the instruction. We design a small comprehension
166 quiz with 7 multiple choice questions and 2 chances to submit to ensure the participants understand
167 the instructions correctly.

168 In the CEG annotation tasks, the participants are asked to label the correctness and causal relation of
169 two event descriptions as shown in Fig. 7. We give 4 examples with detailed explanations to help
170 them understand the rationale of the task. We also give an illustration of object colors and shapes for
171 reference. Bounding boxes are added to the videos to accelerate the process of locating objects of
172 interest.

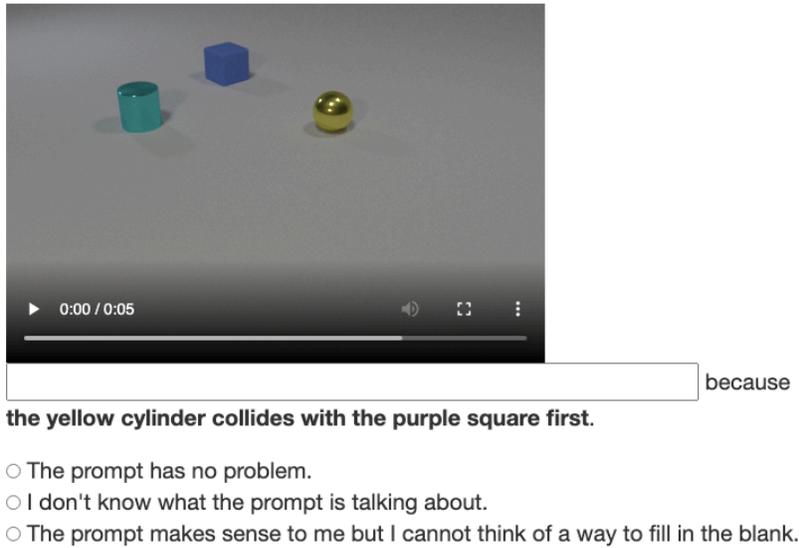
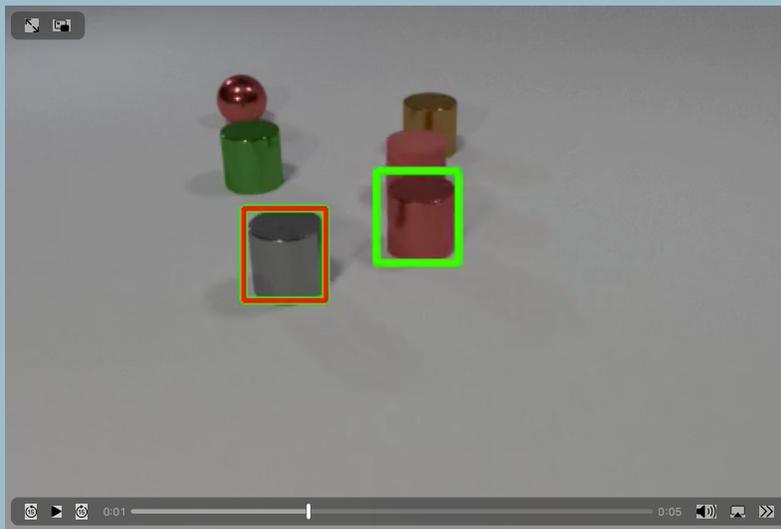


Figure 6: Example of a causal cloze trial. The participants are asked to fill in the blank after watching a video. They can also select the checkboxes if they do not understand the prompt.



Event A: the red cylinder slides into the grey cylinder.

Event B: the grey cylinder moves left.

Question 1: Choose one. Here "incorrect" means either grammatically or factually.

- Event A is incorrect
- Event B is incorrect
- Both event A and B are incorrect
- Both event A and B are correct

Explanation: Choice 4 is selected because both event A and B are grammatically correct and they actually happened in the video.

Question 2: How much is **event A** responsible for **event B**?

- 1 - not responsible at all
- 2 - a bit responsible
- 3 - moderately responsible
- 4 - quite responsible
- 5 - extremely responsible

Explanation: One may think event A is highly responsible for event B as event A is the direct cause of event B. However, there's no right answer to this question - just select the answer you think is most reasonable.

Figure 7: Example of a CEG trial. The participants are first asked to select if the event descriptions are correct. If both correct, they are asked to label the level of causal relations between the descriptions. For each event pair, we provide bounding boxes for objects involved in the events for better annotation efficiency.

173 **D Dataset Release**

174 Our dataset is under CC0 License. We provide a documentation using data statements for NLP in [9].

175 **Short form data statement** CLEVRER-Humans is a large-scale video reasoning dataset of human-
176 annotated physical event descriptions and their causal relations. It contains machine-generated texts
177 based on crowdsourcing data in US English (en-US). The language quality and causal structure
178 annotations are obtained by watching videos, reading texts, and entering responses on MTurk.

179 The following is the long form data statement of CLEVRER-Humans:

180 **Curation Nationale** CLEVRER-Humans contains descriptions and causal relations of physical
181 events such as an object entering the scene and two objects colliding with each other. The goals in
182 selecting texts were to ensure the interpretability and correctness of the descriptions and to provide
183 a variety of free-form captioning of physical events in videos. We first collected human written
184 event descriptions by causal cloze tasks, then used machine learning models to generate more natural
185 language descriptions based on the curated data. We post-processed the data by grammar checking,
186 object existence checking, and verb re-balancing. Finally, we obtained human annotations on the
187 texts through crowdsourcing: if the labelers annotated the texts interpretable and correct, we ask them
188 to provide a pairwise graded causal judgment of the events.

189 **Language Variety** The event description data for causal cloze tasks were collected on MTurk.
190 Information about which varieties of English are represented is not available, but at least CLEVRER-
191 Humans includes US (en-US) mainstream English.

192 **Speaker Demographic** We used a cascaded generator composed of a rule-based event detector and
193 a neural pairwise generator to generate texts. When the curating training data in causal cloze tasks,
194 we restricted the location of these MTurkers to be in the US. It is expected that most speakers use
195 English as their native language. Estimated demographics of MTurkers may refer to [10].

196 **Annotator Demographic** We hire the MTurkers with the approved HITs of 1000 or higher. We
197 expect the MTurkers to be the general public who are familiar with basic crowdsourcing process.
198 When collecting data, we release the tasks in batches, where each HIT contains 30 QA pairs mostly
199 coming from one or two videos. We perform quality check to unsure annotators have sufficient
200 knowledge of English language. We also answer their questions about the annotation process by
201 emails. It is expected that most speakers use English as their native language. Estimated demographics
202 of annotators may refer to [10].

203 **Speech Situation** The intended audience of the texts is the general public. The texts are all in
204 written form. MTurkers are expected to read the text and watch the video when doing causal cloze
205 and causal labeling tasks. The video is about 5 seconds, which can be played as many times as one
206 wishes.

207 **Text Characteristics** The texts are plain English descriptions of a physical event in a video. A
208 sentence usually contains one or more physical object(s) (i.e. sphere, cylinder, or cube) and the
209 related movements or interactions presented in the video. Ideally, the generated event descriptions
210 can maintain the vocabulary and structural characteristics similar to the training data from causal
211 cloze tasks. The detailed statistics of the text are shown in the Dataset Statistics section.

212 **Recording Quality** N/A

213 **Other** N/A

214 **Provenance Appendix** The videos of CLEVRER-Humans are from the CLEVRER dataset [1].

215 **D.1 Intended Use**

216 CLEVRER-Humans can be used as a benchmark in physical scene understanding and causal reasoning.
217 It evaluates machines ability to understand and analyze physical interventions in a restricted setting.
218 Machines are provided with a short video and expected to answer questions regarding the causes of
219 events in the video.

220 **D.2 Maintenance Plan**

221 We will host our dataset permanently on our project’s website. Users are granted access to the dataset
222 through links on the website. We provide versioning of the dataset and archive backup regularly.

223 **D.3 Quality Check**

224 Quality checks over CEG node correctness are performed by majority voting. Since we have split
225 the annotation of each video to 3 annotators, and they will see overlapping events and annotate their
226 correctness. Checks for edge correctness are performed by including additional “quality checking”
227 questions. Specifically, each annotator will see 3 videos and 10 questions for each video. 1 of the
228 video will be from a small and manually-curated dataset by authors, containing 30 videos. The entire
229 answer set will be accepted if and only if the annotators answers those quality-checking questions
230 correctly (more specifically, have a small divergence with our answer).

231 **References**

- 232 [1] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B Tenenbaum.
233 CLEVRER: Collision events for video representation and reasoning. In *ICLR*, 2020. 4, 11
- 234 [2] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of
235 neural machine translation: Encoder-decoder approaches. *arXiv:1409.1259*, 2014. 6
- 236 [3] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780,
237 1997. 6
- 238 [4] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word
239 representation. In *EMNLP*, 2014. 6
- 240 [5] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6, 7
- 241 [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
242 In *CVPR*, 2016. 6
- 243 [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirec-
244 tional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 7
- 245 [8] David Ding, Felix Hill, Adam Santoro, Malcolm Reynolds, and Matt Botvinick. Attention over learned
246 object embeddings enables complex visual reasoning. [https://github.com/deepmind/deepmind-research/
247 tree/a5522d078413e340a2caa1ab0c86d76d2b7efa40/object_attention_for_reasoning](https://github.com/deepmind/deepmind-research/tree/a5522d078413e340a2caa1ab0c86d76d2b7efa40/object_attention_for_reasoning), 2021. 7
- 248 [9] Emily M. Bender and Batya Friedman. Data Statements for Natural Language Processing: Toward
249 Mitigating System Bias and Enabling Better Science. *TACL*, 6:587–604, 12 2018. 11
- 250 [10] Joel Ross, Lilly Irani, M. Six Silberman, Andrew Zaldivar, and Bill Tomlinson. Who are the crowdworkers?
251 shifting demographics in mechanical turk. In *CHI Extended Abstracts on Human Factors in Computing
252 Systems*, 2010. 11