

Data Collection and Worker Quality Control for Generalized and Contextualized Story Explanations (GLuCOSE)

Nasrin Mostafazadeh & Lori Moon
Elemental Cognition
crowd@elementalcognition.com

1 Introduction

GLuCOSE, Generalized and Contextualized Story Explanations, is a large-scale common sense knowledge acquisition task. Given a short story and a sentence X , that is contained in the story, the GLuCOSE task asks crowd workers to compose mini-theories about the world based on the sentence/story pair. These mini-theories are in the form of semi-structured inference rules. There are 20 rules solicited per sentence, guided by ten dimensions of common sense causal explanation with one mini-theory for each dimension being specific to the story and another mini-theory generalized from the specific mini-theory to other situations in the world.

We collected a total of 337,636 pairs of specific and general rules. Of the rules, 28.4% are given the highest quality rank, ‘3’, 44.9% are given the second highest quality rank, ‘2’, and 26.7% are given the quality rating of ‘1’. All rules collected were done by workers who went through qualification and training for the task and our own results, reported in [ACL citation] used all work quality levels in training and development.

This document describes the task, the data collection process, the work quality management, and analysis of the data collection task. The data download is available at <https://github.com/ElementalCognition/glucose>. For results from using the GLuCOSE rules in model training see [ACL 2020 paper].

2 GLuCOSE Task

The GLuCOSE Task involves collecting common sense reasoning that is necessary for understanding stories. There are many dimensions to the common sense rea-

soning that goes into understanding text, but we restrict the task to five types. The five types are (1) Events that cause or enable other events, (2) basic human drives and motivations that cause or enable other events, (3) locations of items, individuals, or events that cause or enable other events (4) possession of items or traits that cause or enable events and (5) other attributes, such as changes of states, that cause or enable events.

For each of these five dimensions, workers considered a sentence from a story. They were asked whether the sentence was connected to one of the five dimensions, either before or after the time of the target sentence, given the story context.¹ We choose not to randomize the dimension order because the task was already long, and the predictable order of questions was something that helped workers form a strategy for doing the task in a manageable amount of time.²

The task was set up as follows, if the story in the task is:

My daughter was born last week. I have a ton of pictures to show you.
She is so cute. She has a ton of hair. I am in love.

For each sentence in the story, workers are asked whether each of the 5 dimensions apply before or after that sentence. Let us take the second sentence as an example:

I have a ton of pictures to show you.

Workers are asked first about the things that occurred before the sentence. They are asked if there is anything in the story that causes or enables the events in the target sentence. If there is something relevant, then they write a specific statement about it. They are given slots for words to help them reduce the sentences, and the connective, in this case, ‘Causes/Enables’ is supplied:

My daughter was born Causes/Enables I have pictures to show you

Next, the worker is asked to make a general statement about the world based on the specific statement that they wrote, for example:

Someone_A has a baby Causes/Enables Someone_A shows Someone_B
pictures of the baby

¹The stories used in the task were from ROC stories [citation]. We restricted the set of stories used to those with a vocabulary no larger than that expected for an eight year old child.

²The average time it took a worker to complete one HIT was 8 minutes and 2 second, and the median was 5 minutes and 8 seconds. These numbers, however, do not account for the fact that workers often open multiple HITs and complete them sequentially.

Workers had a drop-down menu of expressions such as ‘Someone_A’ to help them generalize. Each pair of rules like this is a specific and general rule.

For the next dimension, workers were asked, about the same sentence, if there were any basic human desires or motivations that caused events the sentence, in this case ‘I have a ton of pictures to show you’. The worker had more prompts for this response. For example, they were offered a drop-down list of verbs, shown in table 1. A worker might well respond, ‘Someone_A feels love Motivates Someone_A shows Someone_B pictures of Someone_C (who Someone_A loves). This process is repeated for locations, (Is there a location that enables the events in the sentence?), for possessions, and for any other attribute.

After all five dimensions are considered, the order is reversed and they are repeated. Dimension six asks if the events of the sentence under consideration cause or enable anything to follow. An example is below:

I have pictures to show you Causes/Enables You see pictures

And a general statement could be:

Someone_A shows Someone_B Something_A Causes/Enables Someone_B sees Something_A

This continues for the other four dimensions in this direction. So each of ten dimensions total are solicited for each sentence, in a story context. Not all dimensions are relevant for every sentence, but if a worker does consider a dimension to be relevant, then the worker must fill out both the specific and the general statement about it.

3 Data Collection

Data was collected through an in-house user interface, hosted by Amazon Mechanical Turk.

3.1 Participants

A pool of 1039 workers qualified for our main task via our in-house Qualifying Exam. Of the workers, 373 contributed to the main task. The worker pool was constrained by Mturk-internal ratings (e.g., worker has a high acceptance rate, worker has done at least 100 HITs on Mturk), country codes, and by scores on our qualification task.³ We did not limit how many times workers could take the test, however,

³Country codes were originally restricted to a standard list containing the US, the UK, New Zealand, Australia, South Africa, and Canada, but we opened it up to many other countries per worker

dimension	task connective	Slot Constraints
dim 1 An event that directly causes or enables X	Causes/Enables	none
dim 2 An emotion or basic human drive that motivates X	Motivates	verb [feels, wants, likes] object [curiosity, independence, competition, honor, approval, power, status, romance, success, friendship, belonging, health, safety, livelihood, happy, stressed, angered, disgusted, sad, surprised, fearful, trusting, love, obedient, amazed, disappointment, regret, worthless, aggression, optimistic]
dim 3 A location state that enables X	Enables	verb [am is are] preposition [above, across from, at, below, far from, in, in front of, inside of, near, next to, on top of, outside of]
dim 4 A possession state that enables X	Enables	verb [possesses]
dim 5 Other property besides location, emotional state, and possessions make X possible	Enables	verb [am, is, are, has, have, want, wants, need, needs]
dim 6 An event that is directly caused or enabled by X	Causes/Enables	none
dim 7 An emotion that is caused by X	Motivates	verb [feels, wants, likes] object [curiosity, independence, competition, honor, approval, power, status, romance, success, friendship, belonging, health, safety, livelihood, happy, stressed, angered, disgusted, sad, surprised, fearful, trusting, love, obedient, amazed, disappointment, regret, worthless, aggression, optimistic]
dim 8 A change of location that X results in	Enables	verb [am is are] preposition [above, across from, at, below, far from, in, in front of, inside of, near, next to, on top of, outside of]
dim 9 A change of possession that X results in	Enables	verb [possesses]
dim 10 Other change in property (besides location, emotional state, or possession) that X results in	Enables	verb [am, is, are, has, have, want, wants, need, needs]

Table 1: This table lists each dimension that workers were asked about for each sentence/story pair. Some dimensions contained drop-down choices.

they could only take it once per release, where each release had 1,500 HITs. In order to do well on the test, users were encouraged to carefully read instructions and the examples embedded in the test itself and referred to a shared document that contained general guidelines for the task.⁴

3.2 Qualification Content

The qualification test contained questions testing expertise in three areas: Identifying correct use of the UI slots for language expressions (see Figure 1), recognizing the right level of generalization (see Figure 2), and identifying causes and effects with proper temporal understanding of the stories (see Figure 3).

For the slot use, prospective workers were given correct slot examples and incorrect slot use examples and asked to choose the appropriate ones. The slots were simplified from part of speech tags to allow for easier understanding and to encourage simplification of the language of the stories.

For generalization understanding, prospective workers were presented with a Specific Statement based on the target story. They then had to choose the best general rule to derive from that statement. This type of question tested their understanding of the variable format (e.g., Someone_A), their understanding of the fact that proper names do not belong in general rules, and their understanding that, if something is too general, it does not make sense as a rule.

For understanding causes and effects, users were presented with sentences in the cause and effect structure, one of which was a valid cause-effect in the story and the others which were not.

In our mid-rounds, we started an additional level of testing for entrance called ‘the warm-up task’ as a way to ease workers into the large-scale task, if they did not get 100% on the qualification test, but got over 70% correct.

The warm-up was a sample of HITs like those seen in the main task, but it had simpler sentences and only seven HITs. We requested that workers submit three of the seven HITs, so that we had an ample sample-size to give feedback and see if they were ready for the large-scale task.

requests, including India, China, Germany, Jamaica, the Dominican Republic, Belize, and Nigeria. Many countries that are not in the standard list are already English speaking, and, furthermore, the qualification test itself would be challenging to pass if someone did not already have a good command of English. Since we were not limited to a particular dialect of English, we attempted to add the country code of nearly 50 countries which are English-speaking, but the Mturk interface limited us to 35 countries.

⁴The shared document is at: <https://docs.google.com/document/d/1W8S7y97G9yoAO5qVCHrBB7CTctVsZTOrLCwlgP0dw-U/edit>

Welcome to the Qualification Test for "Explain a Story" HIT!

Click to Read The Instructions for the "Explain a Story" First!

Please answer the following questions to get qualified for the "Explain a Story" HIT.

Test Question 1

Test Question 2

Test Question 3

Test Question 4

Test Question 5

Test Question 6

Test Question 7

Test Question 8

Test Question 9

Test Question 10

Story:

Enzo and Zoe were running a race. Enzo fell. He hurt his knee. Zoe looked back. Enzo was her friend. **She ran to Enzo.** She helped him up.

Let's call the highlighted sentence X = **She ran to Enzo.**

Query:

Consider the likely emotions and basic human drives of the participants in X. Does any of these states of mind/feelings motivate the participant to do X?

Below is a partial answer to the above query:

Step 1: We have composed the following [Specific Statement](#).

The specific statement in natural language is:
 "Zoe felt worried" Motivates "Zoe ran to Enzo"

Which of the following is the correct use of the slots for expressing the above Specific Statement?

☐

Zoe	felt	worried	Motivates	Zoe	ran	to	Enzo
subject	verb	object1		subject	verb	preposition1	object1

☐

Zoe	felt worried	Motivates	Zoe	ran to Enzo
subject	verb		subject	verb

Next

You will see the submit button when you reach the end of the questions.
 Thanks for your hard work! If you encounter any issues, please contact us.

Figure 1: Qualification question about slot use

Welcome to the Qualification Test for "Explain a Story" HIT!

Click to Read The Instructions for the "Explain a Story" First!

Please answer the following questions to get qualified for the "Explain a Story" HIT.

Test Question 1 Test Question 2 Test Question 3 Test Question 4 Test Question 5 Test Question 6 Test Question 7 Test Question 8 Test Question 9 Test Question 10

Story:

Lewis was running for president of the chess club. He wanted the club to make some changes. **He especially wanted the club to change the meeting location.** He was tired of waiting for his parents to pick him up. Lewis made up that he was going to give all the kids soda at chess club. It worked! He got elected. But he did not change the location or buy anyone soda.

Let's call the highlighted sentence X = **He especially wanted the club to change the meeting location.**

Query:

Consider the events that happen after X (or are likely to happen). Does X directly cause any of them, or simply make it possible (i.e., enable it)?

Below is a partial answer to the above query:

Step 1: We have composed the following **Specific Statement**.

The specific statement in natural language is:
 "He wanted the club to change the meeting location" Causes/Enables "He..."

For Step 1, which of the following is the best Specific Statement in terms of content?

☐ He wanted the club to change the meeting location Causes/Enables He got elected.
 subject verb object1 object2 subject verb object1

☐ He wanted the club to change the meeting location Causes/Enables He ran for president of the chess club.
 subject verb object1 object2 subject verb preposition2 object1

Next

You will see the submit button when you reach the end of the questions.
 Thanks for your hard work! If you encounter any issues, please contact us.

Figure 3: Sample qualification question about what events follow

3.3 Main UI Content:

Qualifying workers were able to access large batches of data with no limit on how many HITs they could complete. The main UI displayed a page like in figure 5 for each of the 10 GLUCOSE dimensions in the same order each time. When they said that the dimension was relevant for the selected sentence (by clicking 'Yes..' in the "Your Answer:" box), they were taken to a screen where they could input answers, as shown in figure 5.

For the Specific Statements, users freely entered text within the constraints of some part-of-speech guidance for some dimensions, but, for others, there were constraints on verbs or emotions, as shown in table 1. All preposition slots contained a drop-down list of English prepositions.

For General Rules, the subject position had a drop-down menu of variables for people, places, and things, such as "Someone_A". The verb slot was constrained in the same way as with the corresponding Specific Statement for that dimension.

Read Instructions

[Frequently Asked Questions \(FAQ\)](#)

Please answer the following queries about the story below.

The most thorough and accurate submissions will receive bonuses! We have many more HITs coming.

Query 4

Query 5

Query 6

Query 7

Story:

Jennifer has a big exam tomorrow. She wants to nail the exam. She pulls an all-nighter. **The next day, she is very tired.** Her teacher tells the students that the test is postponed. Jennifer is quite relieved.

Let's call the highlighted sentence X = **The next day, she is very tired.**

An event that directly causes or enables X

Consider the events that happen before X (or are likely to happen). Does any of them directly cause X, or simply make X possible (i.e., enable X)?

Whenever possible, you are encouraged to find the answer from the other sentences in the story. Remember, there are often no right or wrong answers; just give us your intuition.

Click to see an example answer for Query 1

Click to see an example answer for Query 1

Click to see an example answer for Query 1

Your Answer:

☐ No, I can't think of anything really/the query is not applicable to this sentence!

☐ Yes, below is my two-step answer.

You will see the submit button when you reach the end of the queries.

Thanks for your hard work! If you encounter any issues, please contact us.

Figure 4: Main UI

9

Your Answer:

☐ No, I can't think of anything really/the query is not applicable to this sentence!
☒ Yes, below is my two-step answer.

Step 1: Fill in the following blank slots to compose a **Specific Statement.**

Note: the content for X goes in the **highlighted** part.

subject	verb	preposition1	object1	preposition2	object2
---------	------	--------------	---------	--------------	---------

Causes/Enables

subject	verb	preposition1	object1	preposition2	object2
---------	------	--------------	---------	--------------	---------

This specific statement turns into the following natural language form:
 "_____" Causes/Enables "_____"

Step 2: Fill in the following blank slots to compose a **General Rule.**

ⓘ The content for X goes in the **highlighted** part.
 ⓘ This General Rule you write should be a reasonably sound and meaningful rule when read in isolation.
 ⚠ Do not undergeneralize or overgeneralize. Read [this](#) document to learn where the sweet spot is.
 ⓘ For some of the Subject and Object slots, after selecting a general type from the dropdown, a little textfield appears at the bottom for adding attribute clauses. If you need to further specify your general type, you can write an **attribute clause** to achieve the right level of generality.

Someone_A	verb	preposition1	object1	preposition2	object2
-----------	------	--------------	---------	--------------	---------

SUBJECT

- Someone_A
- Someone_B
- Someone_C
- Someone_A and Someone_B
- Some People_A
- Something_A
- Something_B

Causes/Enables

verb	preposition1	object1	preposition2	object2
------	--------------	---------	--------------	---------

The following is the automatically generated natural language form of your General Rule. If it is not correct, try to edit your answer in the above General Rule slots until the following form reads coherently and makes sense.

_____ Enables "_____"

Did you want to express your general rule in a way that was truly not feasible using the slots above? Just write it in the text field below.

Figure 5: When “Yes” is selected for “Your Answer:” on the main UI, workers can input answers in slots

3.4 Data Quality Control Pipeline

For work contributed through the main UI, data quality was controlled through daily monitoring of a percentage of incoming submissions and statistics on average dimensions filled out and time. The percentage of answers reviewed by hand were used to modify worker ratings. Figure 6 shows the strategic flow of worker ratings. Workers enter the task with a score of “-1” then advance to “2” as they become more proficient, getting a bonus increase. The top numeric rating is “3”, which has an additional bonus increase. Select workers with a “3” rating were also moved into “top rated” batches that paid more per HIT and included higher bonuses and incentives. If work quality dropped, workers’ ratings were reduced. If their work was at a risk of degrading corpus quality, they were given a “0” and disqualified from the task. All incoming submissions were approved but did not receive a bonus. Several workers were disqualified and then worked to re-qualify and became top rated workers. The General Guidelines and Qualifying Exam, combined with individual feedback on answers, often provided sufficient remedial learning to get good responses again. Most data quality issues were due to workers trying to rush at the task and not reading the General Guidelines carefully. Figure 7 shows a state of the percentage of contributions to the data by worker score.

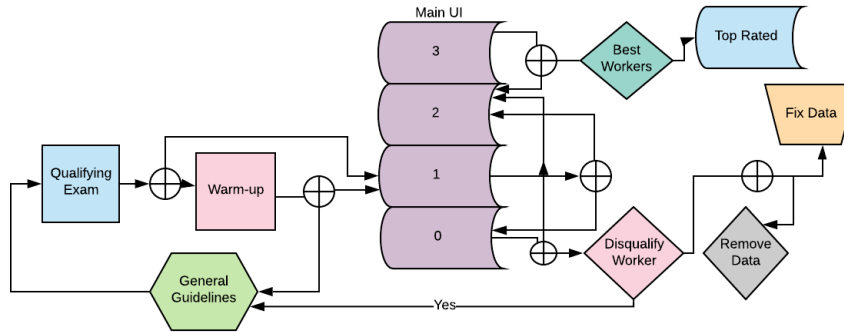


Figure 6: Data collection pipeline. Numeric values in the Main-UI (lavender) correspond to worker ratings.

3.4.1 Review Dashboard

Incoming submissions to the large batches were monitored daily through the in-house review dashboard for quality control. The dashboard is shown in Figure 9.

Percentage of workers with each quality score

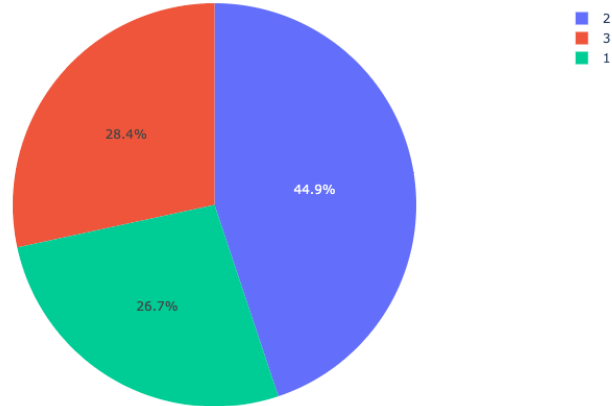


Figure 7: Workers by quality score rating, Red = ‘3’(high), Purple = ‘2’(mid), Green = ‘1’(lower)

Submissions displayed through the review UI could be filtered by batch, submission status (e.g., Submitted or Approved), worker identification number, date of submission, story identification number, dimension, and sentence. When work in need of approval was displayed, it showed the worker’s current quality rating and any notes on the worker’s history.

These ratings were maintained in a spreadsheet that was uploaded to the UI as changes were made.

4 Preparation and Analysis of the Final Data Set

After we completed the collection task, we did additional post-processing of the data. This section discusses that process and the results of the final data.

4.1 Data Post-Processing

In order to get more granularity in the final worker quality ratings, we coordinated the commit date of worker quality rating changes with the submission date of the HITs. This increased granularity improved our quality level to showing 80% very high quality answers when filtering a random 50 samples from workers rated as 3.

Quick MTurk Review Dashboard

Account Balance: 6668.85

Number of all workers qualified for regular

1164

Number of all workers qualified for TopGrade

20

Recount the number of all workers qualified for regular largescale run

Recount the number of all workers qualified for TOP-GRADE largescale run

Worker ID to Disqualify:

Worker ID

Qualification Type ID to disqualify from:

3GR2F62BN49LR89U21N3JKD87TN

Disqualify this worker

Worker ID to QUALIFY:

Worker ID

Qualification Type ID to QUALIFY to:

36F71BC7T3QE48V9ZTN1NRQ84N8

QUALIFY this worker

HIT Type ID:

3JHC69GPMZ4K4MU8XHEG9LJ

Any NextToken to Start with?

Worker ID

Maximum number of assignments to load

2500

Retrieve all the submissions!

Update worker record file

Now you can perform the following actions as you wish!

Qualification Type ID Warmup:

3M0XATTAE9Y0EANZ5V4HCFNSZS

Qualification Type ID Largescale Regular:

3GR2F62BN49LR89U21N3JKD87TN

Qualification Type ID Largescale TopGrade:

36F71BC7T3QE48V9ZTN1NRQ84N8

Review All Qualification Task Submissions and Assign to the Qualification Type ID

Auto-grade ALL submissions!

Dump all submissions to a csv file!

You can filter the list of submissions using the following fields:

HIT ID:

HIT ID

Assignment ID:

Assignment ID

Assignment Submission Time:

Assignment Submission Time

Worker ID:

Worker ID

Story ID:

Story ID

Assignment Status:

Assignment Status

Selected Sentence:

Selected Sentence

Dimensions List:

1,2,3,4,5,6,7,8,9,10

Task Type?

Glucose Main HIT

Loading the next page of the HITs... Sit tight until it finishes. Loaded 1700 hits so far!

Figure 9: Data review dashboard allows reviewers to see user responses, sorted by desired parameters, and give feedback to workers

4.2 Data Quality Results Analysis

In the total data, workers filled in an average of 4.7 dimensions, with the median being 4 dimensions.

Figure 10 shows the distribution of dimensions filled out by type. We expected dimensions 5 and 10 to be less frequent (Is there anything else that causes or enables the sentence?), due to the fact that they refer to rare changes of state, such as becoming wet. We saw some bias towards answering more of the first questions than the later ones, offset by an overall preference for dimensions 1 and 6, the causality dimensions.

Dimensions 8 and 9 were less frequent than their counterparts 3 and 4. The main reason is because a location state is often part of the scene setting or background of a story, whereas a location state that is the result of the highlighted

Percentage of each GLUCOSE dimension in the ENTIRE data

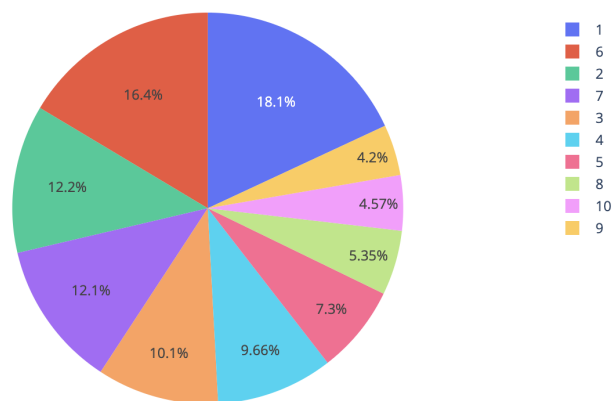


Figure 10: Data collected by dimension in the total data

Percentage of each GLUCOSE dimension in the TEST data

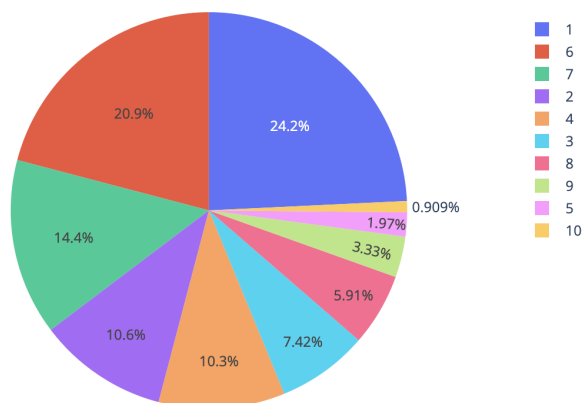


Figure 11: Data collected by dimension in the test data

sentence only occurs when a change of location occurs. Similarly, there are things possessed at the outset of a story, as part of the scene it occurs in, that enable things that are talked about in the story, however, for a sentence to lead to a possession state, it needs to be the case that something is acquired.

Of the final data, we used a portion of the data ranked as quality level ‘3’ to construct a test set. This test set was judged by a group of human annotators and determined to be ‘correct’ or ‘nearly correct’ for all dimensions (see [ACL citation]). We consider the test data to be representative of the larger data set in terms of which dimensions are relevant. Figure 11 shows the percentage of each dimension filled out in the test data.

5 Conclusions

The GLuCOSE task is a large-scale effort to collect common sense mini-theories from crowd workers. We successfully collected rules that prove to be useful in applications such as model training ([cite ACL]).

Using Amazon Mechanical Turk requires breaking tasks down into HITs (Human Intelligence Tasks), and the interface has an economy all its own. There is an expected range of the amount of work per HIT and an expected pay rates for that work. Breaking a large task like GLuCOSE into appropriate tasks required numerous iterations. Having a visually-appealing UI was essential for workers doing this task. Having slots to help with forming mini-theories kept the data in a format that made it easier for applications and provided guidance to workers. Inevitably, this task involved a lot of training for workers and for us as requesters.

6 Acknowledgements

This work could not have been completed without the group of Mturk workers who contributed their time and effort. We are thankful for their patience as we iterated the task instructions and quality control.