

With our flagship DataInFormation<sup>SM</sup> suite of solutions for Image Labeling, Data Annotation, NLP Calibration and Form Data Capture, Liberty Source is proud to serve a wide range of clients across many industry verticals. Our Solution Briefs provide a snapshot of the challenges clients bring to us, and the outcomes we achieve for them.

### Cause & Effect Question Set Evaluation for LLMs

Industry	Prompt Engineering
Project Duration	~2 Months

Data Type	Text
Ongoing?	Yes

### Challenge

Identify potential problems (such as Causal Reasoning and Model Performance) and their causes in our client's existing LLM. Additionally, we needed to verify the performance and enhance the cause/effect relationship in outputs from our client's LLM by generating benchmark "Cause and Effect" question sets.

### Solution

To execute this, our prompt writers began by choosing a word from a list (provided by the client) to use in a sentence. Operating within a set of client parameters, our skilled writers created two statements. Both statements had plausible causes or effects related to the word selected from the original list. Then the writer chose one of the responses to their two statements as the single most accurate response to the client's "Cause" or "Effect" question, thus providing a usable future benchmark.

Those responses were then cycled through an evaluation phase which was based on answering pre-determined questions about the relevance of the content in the response. The post-evaluation responses were then used to teach the model to generate relevant responses in potentially ambiguous situations.

This type of evaluation cannot be done without Human-in-the-Loop (HitL) judgement and expertise, which we utilized by leveraging the domain experience and knowledge of our prompt writers.

### Outcome

We created hundreds of prompts and validated responses (in the form of answer statements) for given "Cause" and "Effect" question prompts. Then, the prompt and answer pairs were exported as a JSON file and fed into a training pipeline for the client's LLM. Lastly, the outcome statements were then used by the client as human-generated benchmarks for training an LLM to accurately determine "Cause and Effect" relationships.

By using these updated benchmark statements, our client successfully pinpointed potential issues and their underlying causes. As a result, they proactively addressed these issues before any problems could arise.

**Input:** List of words or short phrases provided by client, along with direction regarding either cause or effect

Include at least one of these words in the prompt

- goods layer
- , boxes
- , tall tales
- , electric signs
- , door windows



**Process:** HitL generates two effects or two causes and decides which is better - using their domain expertise and based on criteria given by the client. The generated responses are then evaluated to answer a series of questions, resulting in a benchmark dataset for LLM Training.

Include at least one of these words in the prompt

- goods layer
- , boxes
- , tall tales
- , electric signs
- , door windows

Question statement \*

The boxes were left out in the rain.

Question type

What is the effect of this?

Answer 1 \*

The boxes caught on fire.

Answer 2 \*

The boxes fell apart when they were picked up.

Best answer \*

Answer 2

**Output:** HitL-generated cause & effect pairs are then used to train or benchmark LLMs

```
1 {
2   "terms": [
3     "goods layer",
4     "boxes",
5     "tall tales",
6     "electric signs",
7     "door windows"
8   ],
9   "question_statement": "The boxes were left out in the rain.",
10  "question_type": "What is the effect of this?",
11  "answers": {
12    "answer_1": "The boxes caught on fire.",
13    "answer_2": "The boxes fell apart when they were picked up."
14  },
15  "best_answer": "Answer 2"
16 }
```