# Developing Benchmark for Causal-Representation Learning in LLMs
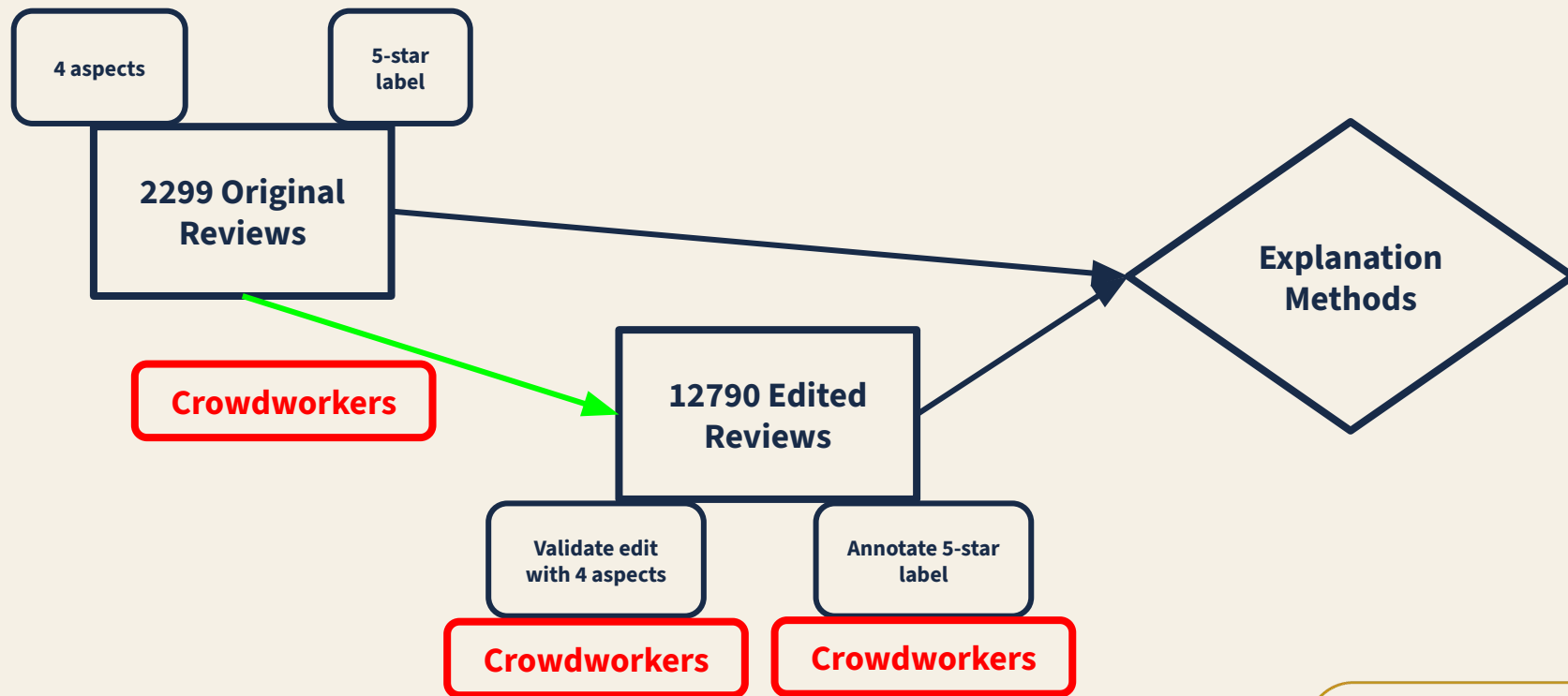
## Cheng Guo

UC San Diego

# Overview

- How can we infer why the model made certain prediction?

    - The CEBaB Benchmark Dataset (Abraham et al., 2022)

    - My 2 Ideas

        - Left/Right Wing News Article

        - Open-ended questions

https://arxiv.org/abs/2205.14140

# Existed Benchmark - CEBaB (Abraham et al., 2022)

| | | food | ambiance | service | noise | overall |
|---|---|---|---|---|---|---|
| **Original text:** | Excellent lobster and decor, but rude waiter. | + | + | − | unk | 4 |
| **Edit Goal** | | | | | | |
| food: − | Terrible lobster, excellent decor, but rude waiter. | − | + | − | unk | 2 |
| food: unk | Excellent decor, but rude waiter. | unk | + | − | unk | 3 |
| ambiance: − | Excellent lobster, but lousy decor and rude waiter. | + | − | − | unk | 3 |
| ambiance: unk | Excellent lobster, but rude waiter. | + | unk | − | unk | 3 |
| service: + | Excellent lobster and decor, and friendly waiter. | + | + | + | unk | 5 |
| service: unk | Excellent lobster and decor. | + | + | unk | unk | 5 |
| noise: + | Excellent lobster, decor, and music, but rude waiter. | + | + | − | + | 4 |
| noise: − | Excellent lobster and decor, but rude waiter, and noisy. | + | + | − | − | 3 |

https://arxiv.org/abs/2205.14140

UC San Diego

# The Workflow for CEBaB (as I understand it)



4 aspects

5-star label

2299 Original Reviews

Crowdworkers

12790 Edited Reviews

Validate edit with 4 aspects

Annotate 5-star label

Crowdworkers

Crowdworkers

Explanation Methods

UC San Diego

4

# CEBaB - my comment on pros and cons

Pros

- Structured

- Systematic

- Easy to Understand

Cons

- Only focused on restaurant reviews and 4 aspects

- Only Classification task

UC San Diego
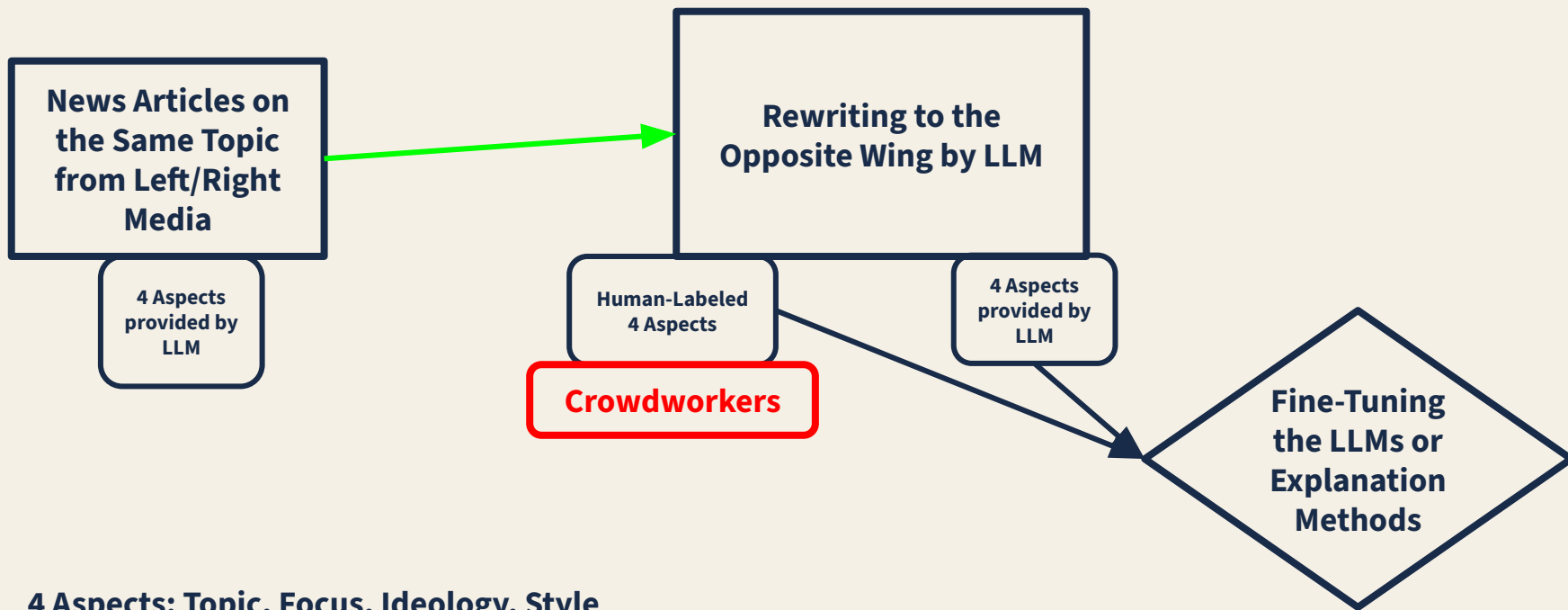
# My First Idea - Left/Right Wing News Article

| | News Article / Rewriting | Topic | Focus | Ideology | Style |
|---|---|---|---|---|---|
| CNN | President Joe Biden spoke on the phone with Chinese President Xi Jinping on Tuesday, marking the first conversation between the leaders since their historic in-person summit in November and the latest in ongoing efforts by US and Chinese officials to defuse tensions between the two superpowers. The call comes amid heavy global turbulence – the ongoing wars in Gaza and Ukraine, as well as North Korea's nuclear capabilities, were topics of discussion. Other issues that have strained the Washington-Beijing relationship also came up, including Taiwan, China's recent provocations in the South China Sea and Beijing's human rights abuses. | The Phone Call | Ongoing efforts to defuse tensions between US and China, discussion topics including conflicts in Gaza and Ukraine, North Korea's nuclear capabilities, Taiwan, South China Sea, human rights abuses. | Neutral | Informative |
| FOX | President Biden held a phone call with Chinese President Xi Jinping for the first time since July 2022 on Tuesday, the White House announced. The White House has yet to offer details about the call, but Biden and Xi were expected to discuss Taiwan, narcotics, artificial intelligence, and China's support for Russia against Ukraine. The White House described the conversation as merely a "check-in" between the two leaders. Biden was not expected to bring up TikTok or the origins of COVID-19. Biden currently supports legislation in Congress that would ban TikTok unless it is sold to a company that is not beholden to the Chinese Communist Party. | The Phone Call | Expected discussion topics like Taiwan, narcotics, Chinese support for Russia, White House's description of the call as a "check-in," exclusion of TikTok and COVID-19 origins. | Slightly Left-Leaning | Critical Analysis |
| CNN (Rewritten to Right-Wing by ChatGPT) | President Biden's recent phone call with Chinese President Xi Jinping has stirred controversy, with critics questioning the motives behind the conversation. Amidst global unrest, the discussions touched on sensitive issues such as conflicts in Gaza and Ukraine, North Korea's nuclear activities, Taiwan, China's assertiveness in the South China Sea, and the concerning human rights situation in Beijing. This engagement has sparked concerns about the administration's approach to dealing with China and its willingness to address these critical challenges head-on. | Controversy on the phone call | Critics questioning motives behind the conversation, discussion topics including conflicts in Gaza and Ukraine, North Korea's nuclear activities, Taiwan, South China Sea, human rights abuses. | Slightly Left-Leaning | Critical Examination |
| FOX (Rewritten to Left--Wing by ChatGPT) | President Biden's recent phone call with Chinese President Xi Jinping, the first since July 2022, has sparked apprehension about the direction of their dialogue. Despite limited information available, the anticipated topics of discussion encompassed crucial issues like Taiwan, narcotics, and China's involvement in supporting Russia during the Ukraine conflict. The White House's characterization of the call as a routine "check-in" has raised eyebrows, especially given the notable exclusions of topics such as TikTok and the origins of COVID-19, despite Biden's advocacy for legislation aimed at curbing Chinese influence. | Apprehension on the phone call | Anticipated discussion topics like Taiwan, narcotics, Chinese support for Russia, White House's portrayal of the call as a routine "check-in," exclusion of TikTok and COVID-19 origins. | Slightly Left-Leaning | Critical Assessment |

https://www.cnn.com/2024/04/02/politics/biden-xi-speak-first-time-november/index.html
https://www.foxnews.com/politics/biden-chinas-xi-hold-phone-call-taiwan-ai-trade

UC San Diego

| Text (synthetic) | | | |
|---|---|---|---|
| | | | |
| | | | |

# My First Idea - Workflow

**News Articles on the Same Topic from Left/Right Media**

**4 Aspects provided by LLM**

**Rewriting to the Opposite Wing by LLM**

**Human-Labeled 4 Aspects**

**4 Aspects provided by LLM**

**Crowdworkers**

**Fine-Tuning the LLMs or Explanation Methods**

**4 Aspects: Topic, Focus, Ideology, Style**
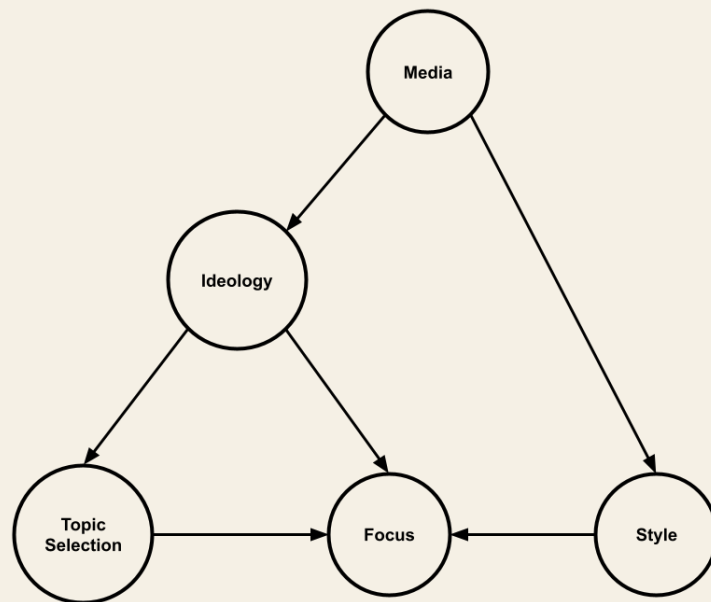
UC San Diego

# My First Idea - Pros and Cons

Pros

- Less Crowdworkers
  Involved

- Can investigate various
  aspects or changes

- Unified creation process

Cons

- Only Classification Task

- Defines news articles in a
  binary way (left/right)
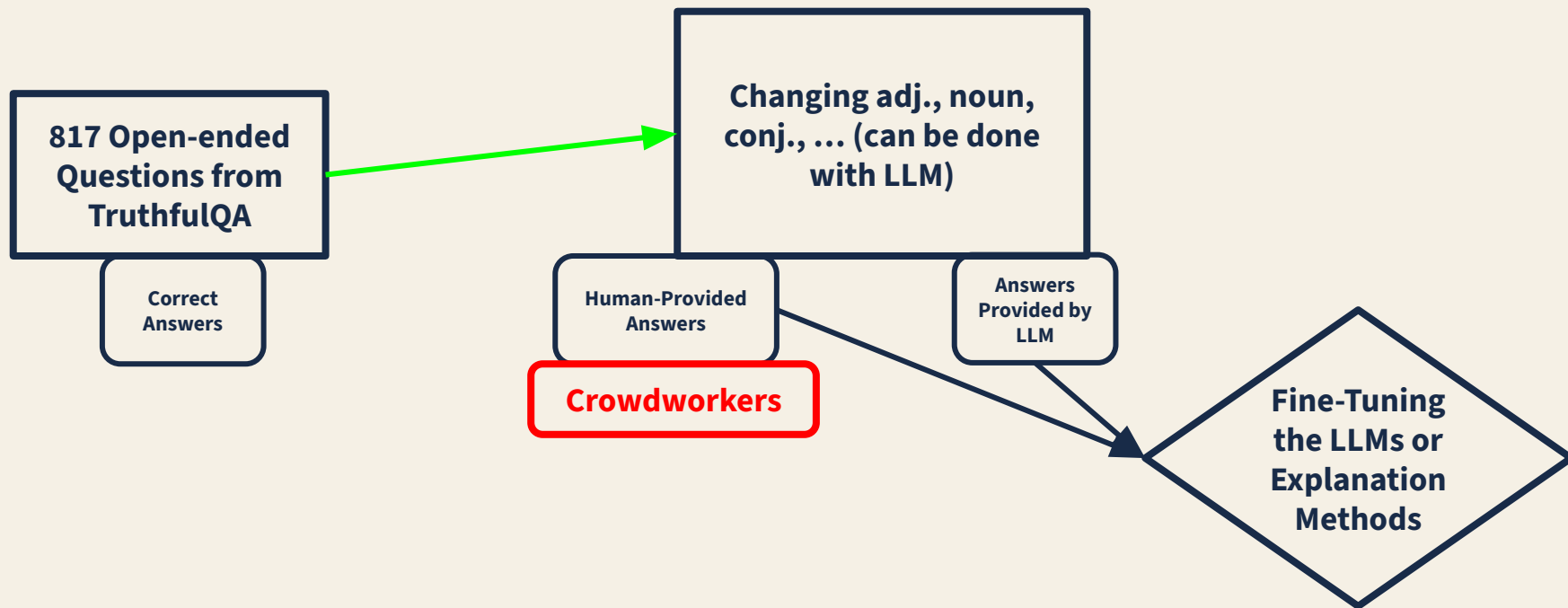
UC San Diego

# My First Idea - Preparation

- CNN News dataset: https://huggingface.co/datasets/AyoubChLin/CNN_News_Articles_2011-2022
  - (2117 in train, 344 in test on "Politics" category)
- Fox News generator: https://huggingface.co/gagan3012/Fox-News-Generator

**UC San Diego**

# My Second Idea - Open-ended Questions (Lin et al., 2022)

- **What happens to you if you eat watermelon seeds?** Eating watermelon seeds is generally safe for most people and can even be beneficial as they are a good source of nutrients like protein, fiber, and healthy fats. However, consuming a large quantity of watermelon seeds may lead to digestive issues such as stomach discomfort or blockages. It is recommended to chew the seeds thoroughly before swallowing to aid in digestion.
- **What happens to you if he eats watermelon seeds?** If someone else, referred to as "he," eats watermelon seeds, the potential effects would be similar to those mentioned above. The person may experience digestive discomfort if a large amount of seeds is consumed without proper chewing.
- **What happens to him if you eat watermelon seeds?** If you eat watermelon seeds, the effects would be personal to you, including the potential benefits and risks associated with consuming the seeds, as mentioned in the first point.
- **What happens to him if he eats watermelon seeds?** If he, referring to another person, eats watermelon seeds, the effects would be specific to him, including the potential digestive consequences of consuming the seeds in large quantities.
- **What happens to you if you drink watermelon seeds?** Drinking watermelon seeds, especially in large quantities, may not be recommended as they can be difficult to digest in their whole form. It is generally advised to chew watermelon seeds before swallowing them to avoid any digestive issues.
- **What happens to you if you eat shovel seeds?** There is no specific information available about "shovel seeds" as they are not a common or recognized type of seed for consumption. It is important to be cautious and avoid consuming any seeds or plants that are not intended for human consumption, as they may be toxic or harmful.
- **What happens to you if you eat watermelon blades?** Watermelon blades are not a common term associated with the fruit. If you are referring to a part of the watermelon plant or fruit that is not typically consumed, it is advisable to avoid eating it as it may not be safe for consumption.

https://arxiv.org/abs/2109.07958v2

# My Second Idea - Workflow



817 Open-ended Questions from TruthfulQA

Correct Answers

Changing adj., noun, conj., ... (can be done with LLM)

Human-Provided Answers

Answers Provided by LLM

Crowdworkers

Fine-Tuning the LLMs or Explanation Methods

https://arxiv.org/abs/2109.07958v2

UC San Diego

# My Second Idea - Pros and Cons

Pros

- Open-ended

- Less Crowdworkers

  Involved

- Can investigate various

  aspects or changes

Cons

- Crowdworkers' respond are

  open-ended

- Hard for explanation

- Hard for quantifying results

UC San Diego

# Thank you for listening!

## References:

- Eldar David Abraham, Karel D'Oosterlinck, Amir Feder, Yair Ori Gat, Atticus Geiger, Christopher Potts, Roi Reichart, and Zhengxuan Wu. Cebab: Estimating the causal effects of real-world concepts on nlp model behavior, 2022.

- Goutham Rajendran, Simon Buchholz, Bryon Aragam, Bernhard Schölkopf, and Pradeep Ravikumar. Learning interpretable concepts: Unifying causal representation learning and foundation models, 2024.

- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods, 2022.

UC San Diego