

# Supplementary material for: WIQA: A dataset for "What if..." reasoning over procedural text

Niket Tandon\*, Bhavana Dalvi Mishra\*, Keisuke Sakaguchi,  
Antoine Bosselut, Peter Clark

Allen Institute for Artificial Intelligence, Seattle, WA  
{nikett,bhavanad,keisukes,antoineb,peterc}@allenai.org

## A Topicwise consistency

We study trends in topic-wise accuracy of models as they read more context information. Bert no-para model does not have access to any context or paragraph, except the language model's background knowledge from Wikipedia. By reading the paragraph context Bert with-para model performs much better on certain topics such as Pollination, blood, mountain, evaporation but the impact of reading is much less on topics such as Igneous rocks, plant crops, solar eclipse, DNA replication. Topics such as blood are very popular on Wikipedia and distributed across several very different articles. These topics are harder for BERT as it requires additional paragraph context to understand the question.

topic	BERT (no para)	BERT
igneous rock	0.66	0.64
plant crops	0.61	0.61
solar eclipse	0.43	0.43
frog	0.59	0.62
DNA replication	0.58	0.63
water cycle	0.63	0.69
fish	0.5	0.57
pumpkin	0.61	0.69
pollination	0.62	0.75
blood	0.62	0.76
mountain	0.57	0.72
evaporation	0.42	0.67

Table 1: As the Bert model (that has access to the paragraph in context) reads more paragraphs in context, its accuracy is better. Reading helps certain topics such as Pollination, blood, mountain, evaporation more than others

\*Niket Tandon and Bhavana Dalvi Mishra contributed equally to this work.

## B Crowdsourcing Influence Graphs

We crowdsource influence graphs by getting the graphs constructed progressively, with the help of five questions stated in Figure 3. At first, the turkers see an empty graph in Figure 1.

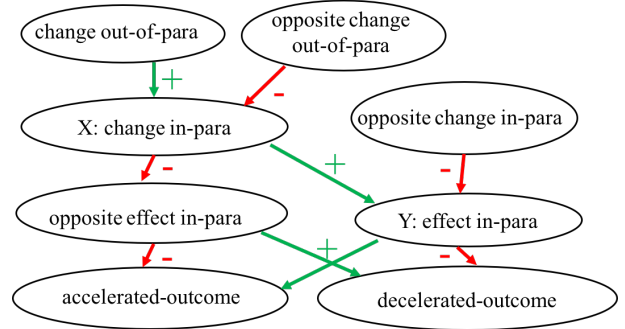


Figure 1: At the start of the process to annotate an influence graph for a given paragraph, the annotators see a blank influence graph with the basic structure.

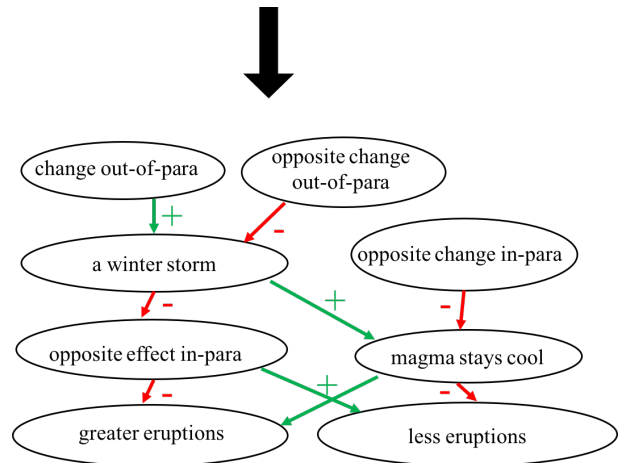


Figure 2: As the annotators answer questions in Fig. 3, a partial influence graph emerges. As they answer questions, the annotators found it useful to validate their answers by examining the emerged influence graph.

When the annotators answer the first question

Consider the paragraph explaining step by step "What causes a volcano to erupt?":

- Magma builds up.
- The magma becomes larger and larger.
- It gets to much.
- The pressure builds.
- The volcano erupts.

1) In this paragraph, suppose this change occurs: . It will have the intermediate effect: , finally resulting in .

2) If the opposite of a winter storm is true, some intermediate effects will be: .

3) Enter some changes  that will result in the opposite of .

4) Imagine! In general,  are reasonable situations that can result in a winter storm.

5) Imagine! In general,  are reasonable situations that can result in the opposite of a winter storm.

**Helpful advice:** An abbreviated example again (from a different paragraph about rain), to show you the style we're looking for:

1. Suppose [the air is warmer], it causes [more water evaporates] -> [more rain]
2. If the opposite of "the air is warmer" is true (i.e., "the air is cooler"), some intermediate effects are [less water evaporates, fewer clouds grow]
3. Changes [the weather is cooler, there is less sunshine] cause the opposite of "more water to evaporate" (i.e., "less evaporation").
4. Imagine! In general, [El Nino forms, longer sunny days] causes the air is warmer.
5. Imagine! In general, [an Arctic chill, a heavy snowstorm] causes the opposite of the air is warmer (i.e., "air is cooler")

Figure 3: The interface shown to the annotators on Mechanical turk platform. Given a paragraph in yellow background, the annotators answer the five questions and an influence graph emerges from their answers.

(shown in Fig. 3), two nodes of the partial influence graph are filled (depicted in Fig. 2).

Once all the questions are answered, the influence graph will be ready. During the process of annotation, there are appropriate validations for quality control.

## C Sample Influence Graphs

To get an impression of our crowdsourced influence graph repository, we display four paragraphs (not hand picked) in Figures 4, 5, 6, 7. These range from natural process, to human body process and mechanical process.

### 13. Describe the process of evaporation

Water is exposed to heat energy, like sunlight.  
The water temperature is raised above 212 degrees fahrenheit.  
The heat breaks down the molecules in the water.  
These molecules escape from the water.  
The water becomes vapor.  
The vapor evaporates into the atmosphere.

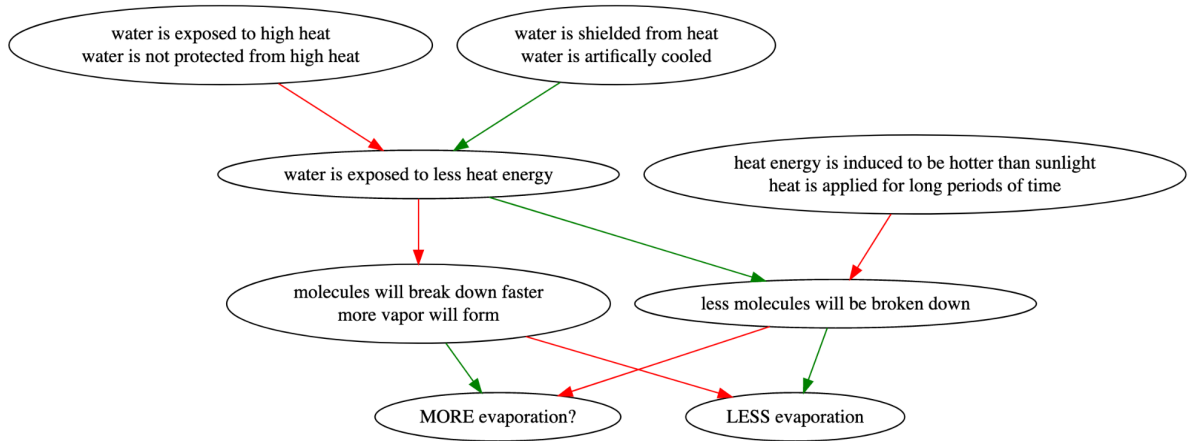


Figure 4: Influence graph for a paragraph from the topic evaporation

### 451. Describe how a flashlight works

Batteries are put in a flashlight.  
The flashlight is turned on.  
Two contact strips touch one another.  
A circuit is completed between the batteries and the lamp.  
The lamp in the flashlight begins to glow.  
The reflector in the flashlight directs the lamps beam.  
A straight beam of light is generated.  
The flashlight is turned off.  
The circuit is broken.  
The beam is no longer visible.

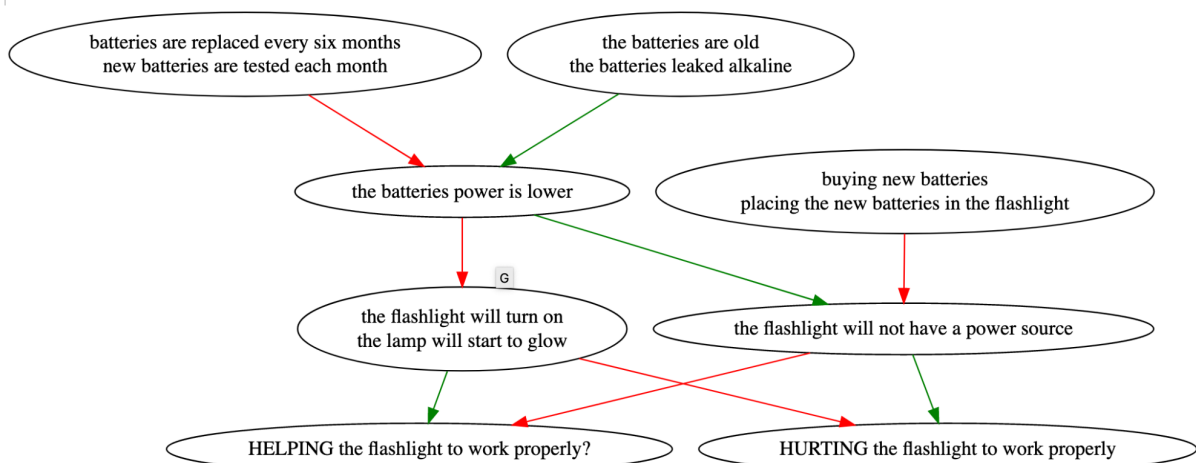


Figure 5: Influence graph for a paragraph from the topic flashlight

#### 16. What do lungs do?

You breathe oxygen into your body through the nose or mouth.  
The oxygen travels to the lungs through the windpipe.  
The air sacs in the lungs send the oxygen into the blood stream.  
The carbon dioxide in the blood stream is transferred to the air sacs.  
The lungs expel through the nose or mouth back into the environment.

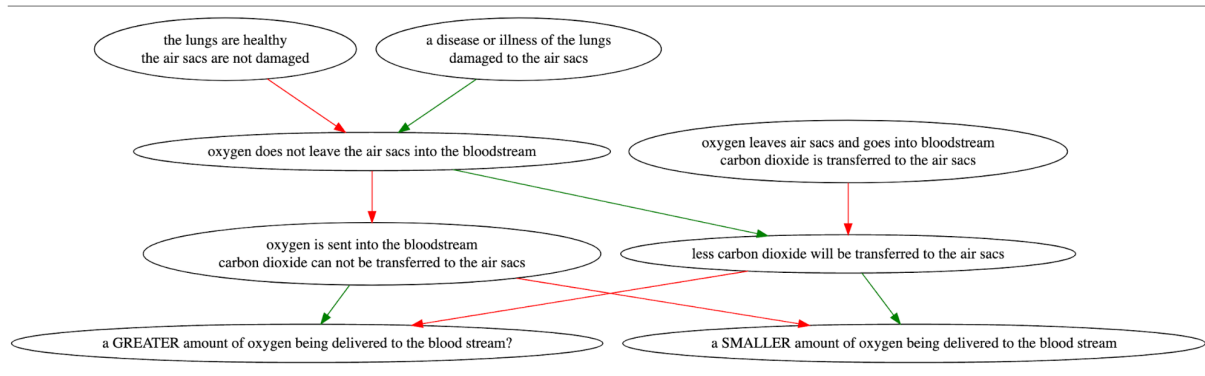


Figure 6: Influence graph for a paragraph from the topic lungs

#### 5. How do minerals form?

Magma comes up to the surface of the earth.  
The magma cools.  
Particles inside the magma move closer together.  
Crystals are formed.  
The crystals contain minerals.

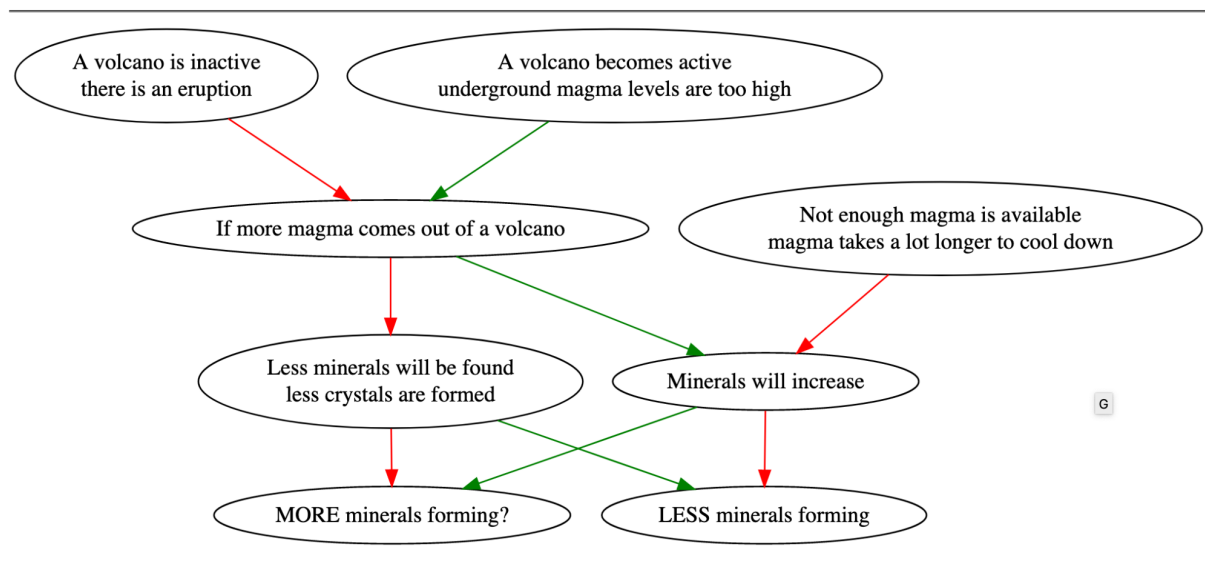


Figure 7: Influence graph for a paragraph from the topic minerals