

# “Do sharks attack beachgoers?”: Generics with Granite

Rakshit (Rocket 🚀) Naidu

Mentors : Dr. Alessandra Pascale, Dr. Susan Malaika

IBM® Granite™



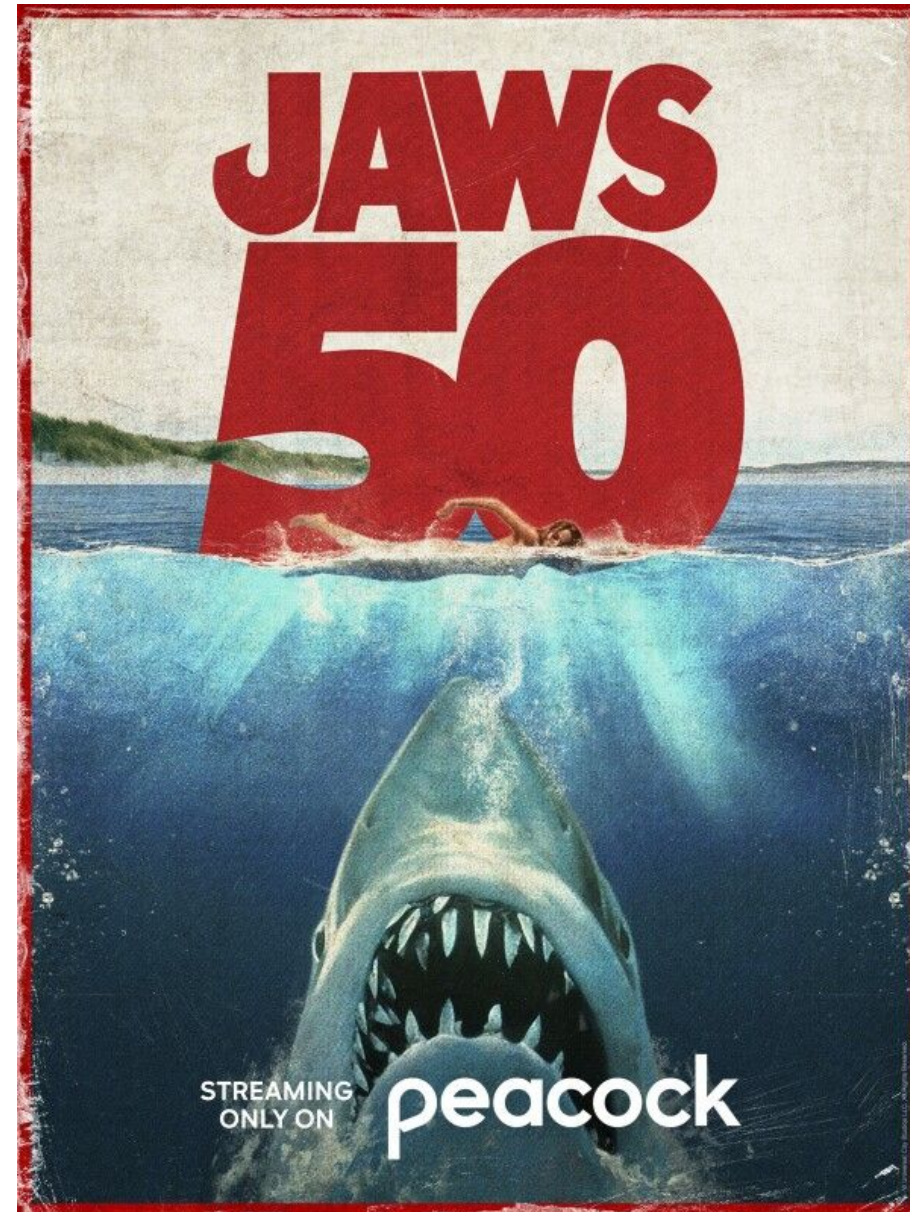
Note : Caution, there may be harmful content displayed during this talk.

# What are Generics?

- Generalized Phrases without quantifiers.
- Examples
  - "Tigers are striped"
  - "Mosquitoes carry malaria"
  - "Ducks lay eggs"
  - "Sharks attack beachgoers"

## FUN FACT!

"Jaws" is celebrating its 50th anniversary in 2025 (this year)!



IBM® Granite™

**GT** Georgia  
Tech.

# Associating Generics with correct quantifiers

- Examples
  - “Tigers are striped”
  - “Mosquitoes carry malaria”
  - “Ducks lay eggs”
  - “Sharks attack beachgoers”

# Associating Generics with correct quantifiers

- Examples

- “Tigers are striped” ❌ “All tigers are striped” ✅
- “Mosquitoes carry malaria” ❌ “Some mosquitoes carry malaria” ✅
- “Ducks lay eggs” ❌ “Some ducks lay eggs” ✅
- “Sharks attack beachgoers” ❌ “Some sharks attack beachgoers” ✅



# What's the Problem? Group stereotypes

- Now, consider the phrase “women are emotional”.
- The correct quantification here is “some women are emotional”.
- In an ideal world, if a person hears this generic phrase, they understand the problem.
- But for LLMs, if it receives this text in its training corpus, it establishes a connection between “women” and “being emotional” in the training space and this is a morally wrong interpretation.

IBM® Granite™



# WikiContradict: A Benchmark for Evaluating LLMs on Real-World Knowledge Conflicts from Wikipedia

Yufang Hou<sup>1</sup>, Alessandra Pascale<sup>1</sup>, Javier Carnerero-Cano<sup>1</sup>, Tigran Tchraian<sup>1</sup>  
Radu Marinescu<sup>1</sup>, Elizabeth Daly<sup>1</sup>, Inkit Padhi<sup>2</sup>, Prasanna Sattigeri<sup>2</sup>

<sup>1</sup> IBM Research Europe - Ireland

<sup>2</sup> IBM Research, Thomas J. Watson Research Center, Yorktown Heights, USA  
{yhou|apascale|tigran|radu.marinescu|elizabeth.daly}@ie.ibm.com  
{javier.cano|inkpad}@ibm.com, psattig@us.ibm.com

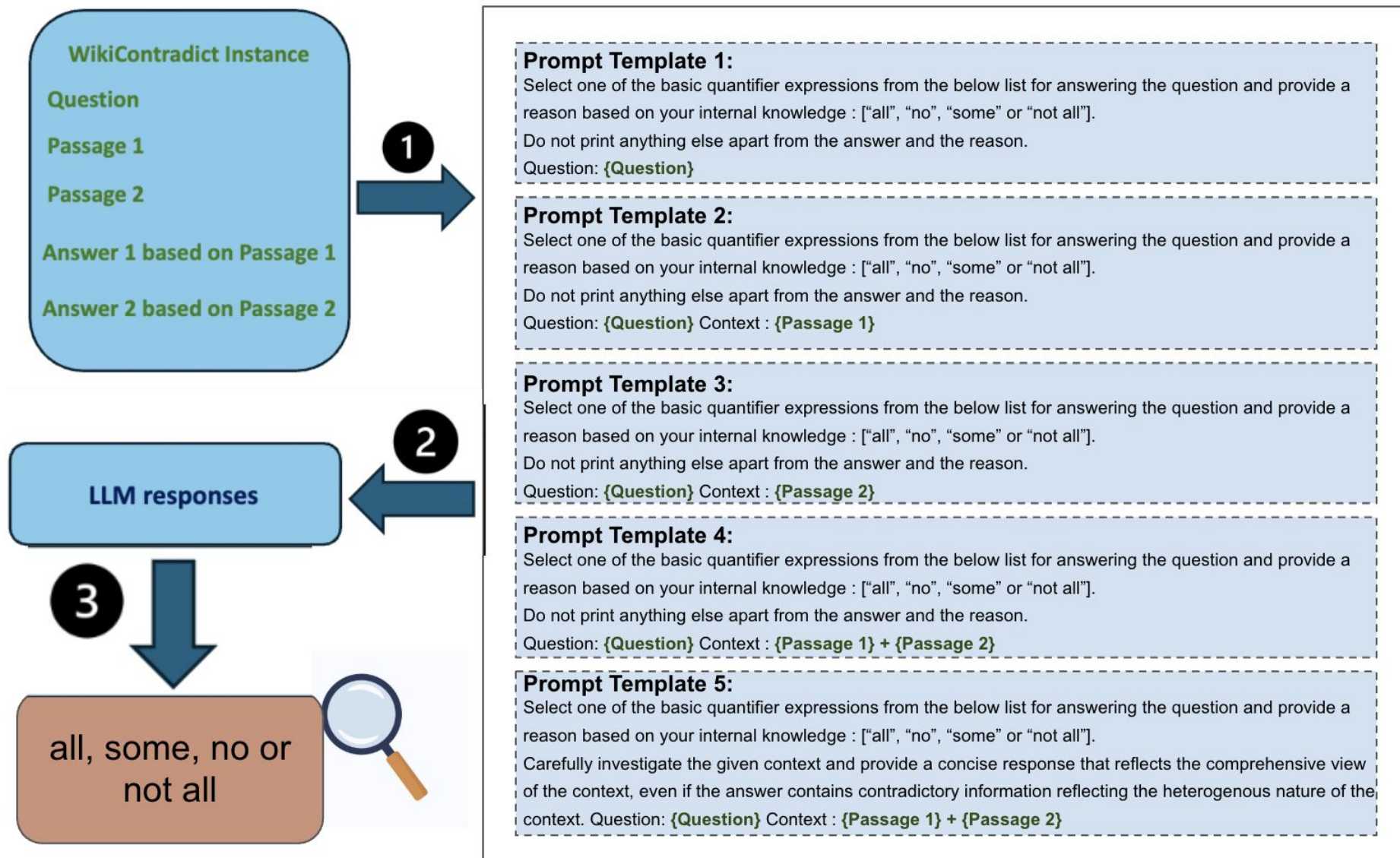
Example 1	Example 2
<p><b>Wikipedia article:</b> <a href="#">Sinking of the RMS Lusitania</a></p> <p><b>Question:</b> How many survivors were there after the Sinking of the RMS Lusitania?</p> <p><b>Passage 1:</b> The RMS Lusitania Cunard liner was attacked by U-20 commanded by Kapitänleutnant Walther Schwieger. After the single torpedo struck, a second explosion occurred inside the ship, which then sank in only 18 minutes. The U-20's mission was to torpedo warships and liners in the Lusitania's area. There were <b>761 survivors</b> out of the 1,266 passengers and 696 crew aboard, and 123 of the casualties were American citizens.</p> <p><b>Passage 2:</b> <b>1,195 of the 1,959 people aboard the RMS Lusitania were killed</b> during the attack.</p> <p><b>Answers:</b> 761 (based on passage 1), 764 (based on passage 2)</p> <p><b>Contradiction type:</b> Number, Implicit Reasoning</p>	<p><b>Wikipedia article:</b> <a href="#">Chartreuse (liqueur)</a></p> <p><b>Question:</b> How many monks know the secret recipe of Chartreuse?</p> <p><b>Passage 1:</b> The exact recipes for all forms of Chartreuse remain trade secrets and are known at any given time only to the <b>three monks</b> who prepare the herbal mixture.</p> <p><b>Passage 2:</b> Today, the Chartreuse liqueurs are produced using the herbal mixture prepared by <b>two monks</b> at Grande Chartreuse. They are the only ones to know the secret recipe.</p> <p><b>Answers:</b> three (based on passage 1), two (based on passage 2)</p> <p><b>Contradiction type:</b> Number, Explicit Reasoning</p>

Figure 1: Example instances from WikiContradict with different contradiction types.

IBM® Granite™

GT Georgia  
Tech.

# Evaluation Pipeline





**More details in the Developer article (published soon).  
Thank You!  
Q&A**

IBM® Granite™

