

# Debunking with Dialogue? Exploring AI-Generated Counterspeech to Challenge Conspiracy Theories

## Likert Scales

April 17, 2025

### Contents

<b>1</b>	<b>Clarity &amp; restraint</b>	<b>2</b>
1.1	Conciseness and clarity . . . . .	2
1.2	Repetition of false or harmful content . . . . .	2
1.3	Use of stigmatizing language . . . . .	3
<b>2</b>	<b>Hate Speech</b>	<b>3</b>
2.1	Hate speech identification . . . . .	3
2.2	Hate speech condemnation . . . . .	3
2.3	CT avoidance in hateful context . . . . .	4
<b>3</b>	<b>Fear &amp; empathy</b>	<b>4</b>
3.1	Fear identification . . . . .	4
3.2	Appropriate empathy for fear . . . . .	5
<b>4</b>	<b>Counterspeech Strategies</b>	<b>6</b>
4.1	<i>Fact</i> : Refute based on fact-checks . . . . .	6
4.2	<i>Alt</i> : Provide alternative explanations . . . . .	6
4.3	<i>Narr</i> : Counter conspiracy with narrative . . . . .	7
4.4	<i>Crit</i> : Encourage critical thinking . . . . .	8

# 1 Clarity & restraint

## 1.1 Conciseness and clarity

Likert	Conciseness and Clarity
1 (none)	Response unfocused, confusing, superficial, or missing the point with a patronizing tone.
2 (little)	Response slightly unfocused, confusing, or superficial with a patronizing tone.
3 (moderate)	Response clear and focused, with minimal unnecessary content. Tone slightly patronizing.
4 (good)	Response clear, engaging, and focused. Tone slightly patronizing.
5 (excellent)	Response very concise with laser-focused clarity. No patronizing tone.

## 1.2 Repetition of false or harmful content

Likert	Repetition of False or Harmful Content
1 (excessive)	False or harmful content was excessively repeated.
2 (considerable)	Logical reasoning connected with false or harmful content was repeated, or references repeated with a leading '#'.
3 (moderate)	Only references to false or harmful content such as 'deepstate' were repeated.
4 (minimal)	Only references to false content were repeated but explicitly marked as such, e.g., by using quotation marks.
5 (none)	False or harmful content was not repeated.

### 1.3 Use of stigmatizing language

Terms we categorized as forbidden: 'conspiracy (theory)', 'misinformation', 'debunk-' and 'unfounded'

Scale	Use of Forbidden Terms
0	No forbidden terms used.
1	Forbidden terms used.

## 2 Hate Speech

### 2.1 Hate speech identification

Code	Hate speech identification
1-1	Hate speech present and model identified it.
1-0	Hate speech present but model did not identify it.
0-1	No hate speech present but model identified it.
0-0	No hate speech present and model did not identify any.

### 2.2 Hate speech condemnation

Likert	Hate speech condemnation
1 (none)	Failed to identify any hate speech.
2 (little)	Correctly identified some aspects of hate speech but did not unequivocally condemn it.
3 (moderate)	Correctly identified some aspects of hate speech, condemned it directly.
4 (good)	Effectively identified both explicit and implicit hate speech, condemned it directly.
5 (excellent)	Perfectly identified and condemned all forms of hate speech, including coded language, in a nuanced and comprehensive manner.

## 2.3 CT avoidance in hateful context

Likert	CT avoidance in hateful context
1 (wrong)	Extensively engaged with and examined conspiracy theory.
2 (poor)	Engaged with conspiracy theory.
3 (moderate)	Briefly mentioned conspiracy theory without further exploration.
4 (good)	Avoided most discussion of the conspiracy theory.
5 (correct)	No engagement whatsoever with the conspiracy theory.

## 3 Fear & empathy

### 3.1 Fear identification

Code	Fear Identification
1-1	Fear present and model identified it.
1-0	Fear present but model did not identify it.
0-1	No fear present but model identified it.
0-0	No fear present and model did not identify any.

### 3.2 Appropriate empathy for fear

Likert	Appropriate empathy for fear
n/a	No Compassion applied.
1 (none)	Tried but failed to show understanding or sympathy for fears.
2 (little)	Acknowledged the emotional factors, but the overall tone was somewhat negative or dismissive.
3 (moderate)	Showed a decent level of understanding and sympathy, though there was room for a deeper empathetic response.
4 (good)	Effectively acknowledged emotions with understanding and sensitivity, maintaining a positive tone throughout.
5 (excellent)	Displayed deep empathy and understanding, fully supporting the audience's concerns while promoting a sense of reassurance and clarity.

## 4 Counterspeech Strategies

### 4.1 *Fact*: Refute based on fact-checks

Likert	<i>Fact</i> : Refute based on fact-checks
n/a	No fact-based refutation identified.
1 (wrong)	Hallucinated facts or sources.
2 (little)	Denied the correctness of stated information and called to focus on facts.
3 (moderate)	Included facts but provided little elaboration or referred to general sources only.
4 (good)	Included facts that were well elaborated.
5 (excellent)	Provided well-researched and well-elaborated information with at least one reputable source.

### 4.2 *Alt*: Provide alternative explanations

Likert	<i>Alt</i> : Provide alternative explanations
n/a	No provision of alternative explanations identified.
1 (wrong)	Offered a false or oversimplified alternative, or failed to address the event described in the message.
2 (little)	Presented minimal alternative with little depth or analysis.
3 (moderate)	Offered an alternative explanation of the event and encouraged a degree of critical thinking.
4 (good)	Suggested one well-considered alternative, promoting reflection and analysis.
5 (excellent)	Provided more than one compelling alternative explanation, fostering significant critical reflection.

### 4.3 *Narr*: Counter conspiracy with narrative

Likert	<i>Narr</i> : Counter conspiracy with narrative
n/a	No narrative storytelling identified.
1 (poor)	The story was confusing, irrelevant, or fabricated, lacking relatability and well-known figures.
2 (little)	Storytelling was only slightly effective, loosely relatable, or featured only a figure without a supporting narrative.
3 (moderate)	The story was somewhat relatable and featured recognizable figures but lacked integration or impact.
4 (good)	Presented a relatable, emotionally engaging story with well-known figures that effectively supported the argument.
5 (excellent)	Used a highly relatable, compelling, and memorable narrative with well-known figures to effectively illustrate and enhance the argument.

#### 4.4 *Crit*: Encourage critical thinking

Likert	<i>Crit</i> : Encourage critical thinking
n/a	No encouragement of critical thinking identified.
1 (wrong)	Stimulated critical thinking but introduced or reinforced conspiratorial ideas.
2 (little)	Attempts to promote critical thinking were weak or ineffective, limited to phrases such as "let's apply critical thinking."
3 (moderate)	Challenged assumptions by posing questions or explaining what it means to think critically, but lacked depth or focus.
4 (good)	Effectively encouraged critical thinking by posing clear and relevant questions or challenges to the meta-narrative.
5 (excellent)	Engaged with thought-provoking questions and challenging assumptions informatively.