

# W

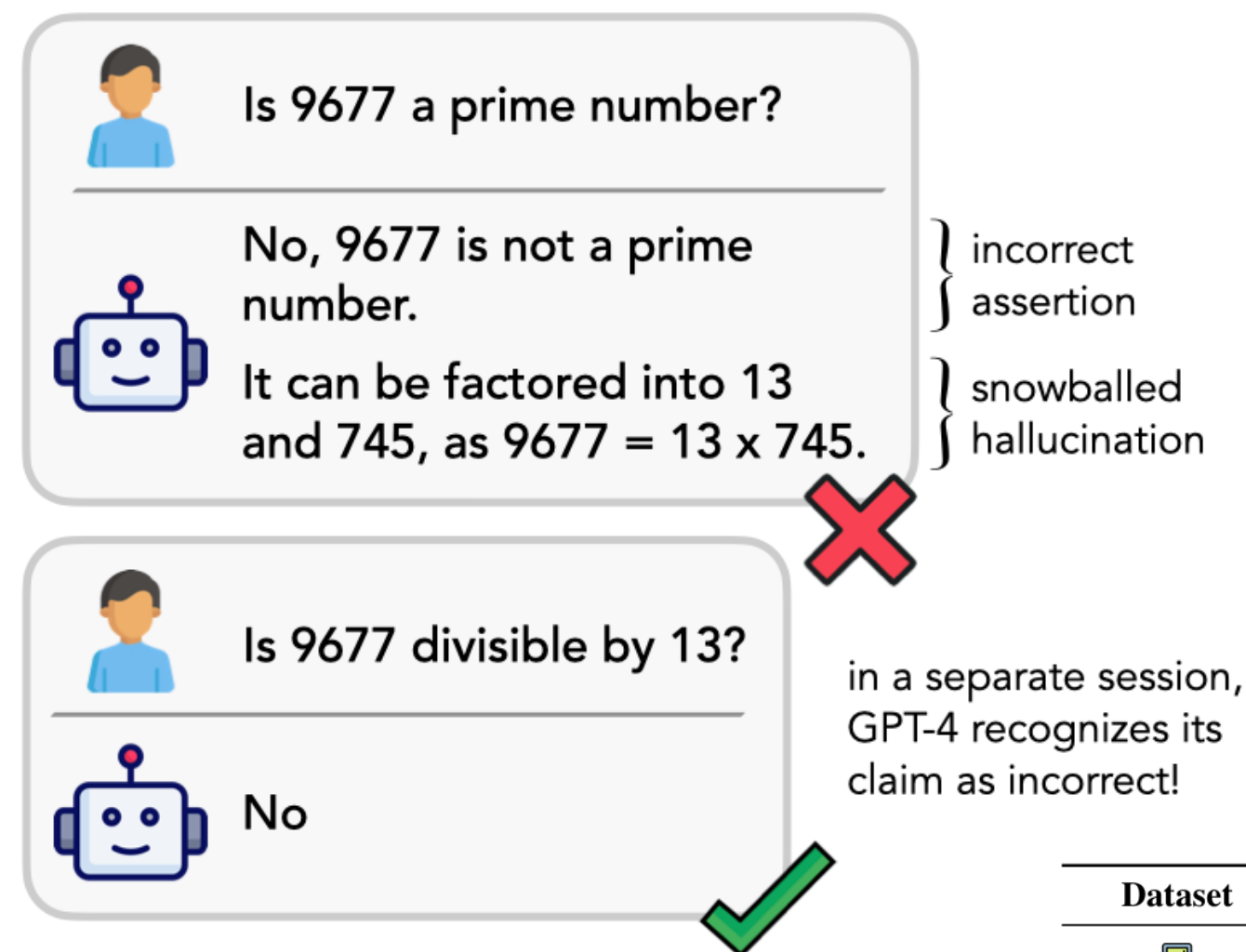
# How Language Model Hallucinations Can Snowball

Muru Zhang<sup>1</sup>, Ofir Press<sup>2</sup>, William Merrill<sup>3,4</sup>, Alisa Liu<sup>1</sup>, Noah A. Smith<sup>1,4</sup>

<sup>1</sup> Paul G. Allen School of Computer Science and Engineering, <sup>2</sup>Princeton University, <sup>3</sup>New York University, <sup>4</sup>Allen Institute for AI

## Overview

- Hallucination is a major issue of language models (LM) where the model gives incorrect claims confidently.
- Hallucinations are often attributed to the knowledge gap in LMs, i.e.: LMs “don’t know” the correct answer.
- In this work, we found **hallucinations can snowball**: previous hallucinations can lead to hallucinations that even the LM itself can reliably identify.
- This highlights LM’s tendency to prefer fluency over factuality.

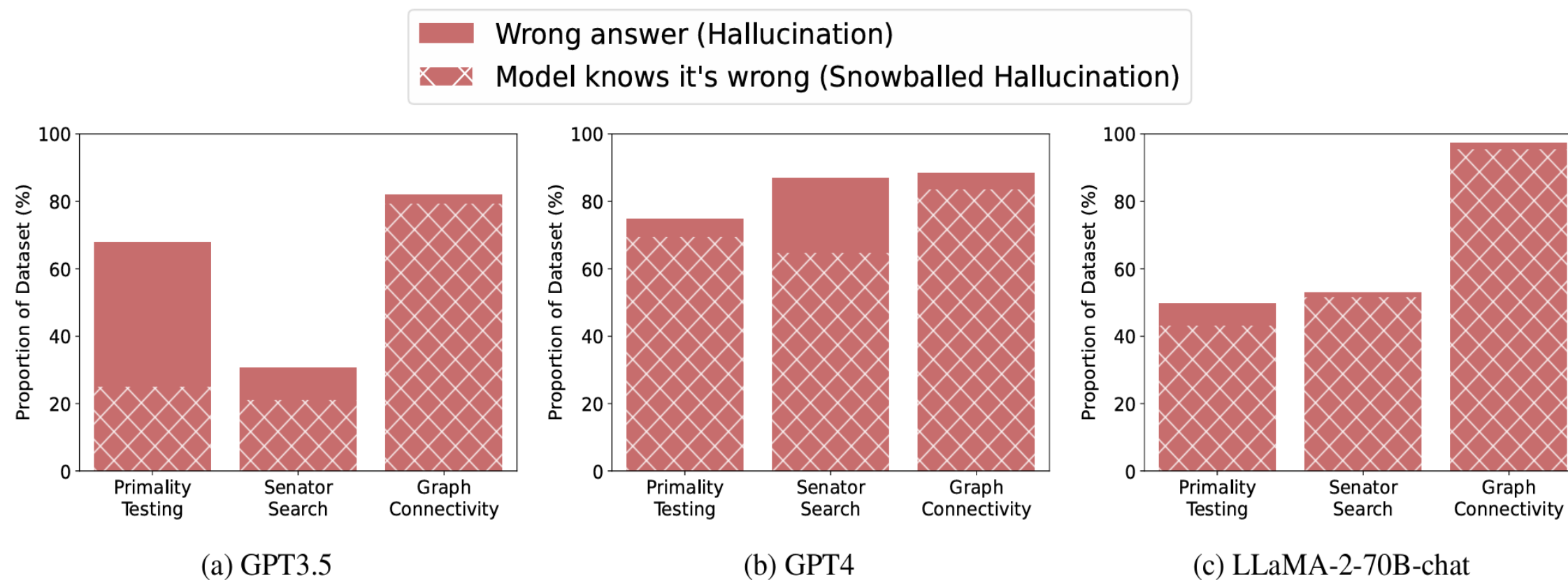


## Experiment Setup

- Three customized QA tasks such that the answer needs to be supported by an evidence.
- 500 questions in each dataset, we define **snowball rate** as the percentage of a LM’s incorrect claims that it can recognize.
- Beam search, higher temperature sampling, zero-shot chain-of-thought prompting all fail to fix the hallucination snowballing issue.

Dataset	Original Question	Verification Question
Primality Testing	<b>User:</b> Is 10733 a prime number? <b>GPT-4:</b> No... It can be factored into $3 \times 3577$ .	<b>User:</b> Is 10733 divisible by 3? Answer with either Yes or No. <b>GPT-4:</b> <u>No</u>
Senator Search	<b>User:</b> Was there ever a US senator that represented the state of New Hampshire and whose alma mater was the University of Pennsylvania? <b>GPT-4:</b> Yes... His name was John P. Hale	<b>User:</b> Was John P. Hale’s alma mater University of Pennsylvania? <b>GPT-4:</b> <u>No</u> ... [it] was Bowdoin
Graph Connectivity	<b>User:</b> Current flight information (the following flights are one-way only, and all the flights available are included below): There is a flight from city F to city K There is a flight from city H to city A [... 10 other rules cut for space ...] Question: Is there a series of flights that goes from city B to city E? <b>GPT-4:</b> Yes... the route is as follows: ... <u>City K to City G</u> ...	<b>User:</b> [...flight information given in the context...] Based on the above flight information, is City K to City G a valid flight? <b>GPT-4:</b> <u>No</u> , based on the above flight information, there is no direct flight from City K to City G.

## Results



## Desired Behavior

- The high recognition accuracy shows that the model have the necessary knowledge to realize the incorrect claim.
- It should refine from giving incorrect claims rather than snowballing from the previous hallucinations.
- Indeed we found this happened on GPT-4 in very small number of cases.