

---

# CSE 584 Final Project: Synthetic Faulty Dataset generation and benchmarking SOTA LLMs

---

Debarshi Kundu<sup>\* 1</sup>

## Abstract

Consider the problem: "If one man and one woman can produce one child in one year, how many children will be produced by one woman and three men in 0.5 years?" Current large language models (LLMs) such as GPT-4o, GPT-o1-preview, and Gemini Flash frequently answer "0.5," which does not make sense. While these models sometimes acknowledge the unrealistic nature of the question, in many cases (8 out of 10 trials), they provide the nonsensical answer of "0.5 child." Additionally, temporal variation has been observed: if an LLM answers correctly once (by recognizing the faulty nature of the question), subsequent responses are more likely to also reflect this understanding. However, this is inconsistent.

These types of questions have motivated us to develop a dataset of science questions where the questions themselves are intentionally faulty. We observed that LLMs often proceed to answer these flawed questions without recognizing their inherent issues, producing results that are logically or scientifically invalid. By analyzing such patterns, we developed a novel method for generating synthetic datasets to evaluate and benchmark the performance of various LLMs in identifying these flawed questions. We have also developed novel approaches to reduce the errors.

## 1. Introduction

### 1.1. Key Contributions

1. **Dataset Creation:** We curated a new dataset of faulty questions, aimed at assessing LLMs' ability to recognize and respond appropriately to flawed questions.

---

<sup>\*</sup>Equal contribution <sup>1</sup>Department of CSE, University of Pennsylvania State University, State College, USA. Correspondence to: Debarshi Kundu <dqk5620@psu.edu>.

Alongside this, a novel technique for generating synthetic datasets was developed. The technique involves the following steps:

- (a) Pick a dataset that contains science questions, such as SciQA and SciQ.
- (b) Extract each question, its corresponding answer, and any additional available information from the dataset for each row item.
- (c) Utilize multiple LLMs, referred to as LLM\_generators (e.g., LLM\_generator1 = GPT 4-o, LLM\_generator2 = Gemini Pro, LLM\_generator3 = Llama 3.1, LLM\_generator4 = Mixtral), to generate faulty versions of the original questions. Each generator also provides a reason why the generated question is faulty and identifies the type of fault, such as logical fallacies, unrealistic scenarios, or violations of physical laws.
- (d) Feed the faulty questions generated in step 3 to another LLM, designated as the LLM\_discriminator (e.g., GPT-4). The LLM\_discriminator is not provided with the reasons for faultiness. Instead, it is tasked with analyzing each faulty question one by one to determine if it is indeed faulty and, if so, to explain why. If the question is not deemed faulty, the LLM\_discriminator answers it. These responses (questions and explanations) are then fed back to the LLM\_generators, which use this feedback to refine and generate a new set of faulty questions.
- (e) Repeat the process: Each LLM\_generator creates one new faulty question per iteration, following the same steps as before, while the LLM\_discriminator evaluates these new questions. This iterative process continues until the LLM\_discriminator can no longer find faults in any of the generated faulty versions of the original question, or until a predefined maximum number of iterations is reached.

Using this iterative approach, we aim to generate faulty versions of the SciQ and SciQA datasets separately. The resulting datasets will include the following columns:

- Science Discipline (and optionally subcategories).
- Original Question.
- Generated Faulty Question.
- Reason why the question is faulty.
- Faulty Answer by the LLM discriminator (applicable only when the discriminator fails to recognize the question as faulty).

2. **LLM Evaluation:** We systematically evaluated the performance of different LLMs, measuring their ability to detect and handle faulty questions. Our findings indicate that current LLMs exhibit varying degrees of expertise across different types of fallacies.

3. **Proposed Error-Reduction Methods:** To address these challenges, we proposed several strategies:

- **AI Agents:** Creating multi-agent systems where multiple LLM models verify each other’s responses before delivering a final answer. This approach leverages the strengths of different models to improve overall performance.
- **Tool Integration:** Incorporating external tools such as WolframAlpha, calculators, fact-checkers, and online search engines into chatbots can significantly enhance their ability to identify and respond appropriately to faulty questions.
- **Harnessing Model Specializations:** Our data suggests that different LLMs have varying areas of expertise. By combining multiple models in a multi-agent framework, we can harness these strengths to create a more robust application capable of effectively addressing flawed questions.
- We also investigated whether certain strategies during training or fine-tuning could make LLMs better at recognizing and addressing such cases. Potential methods include:
  - (a) **Exposure to Faulty Questions:** Introducing flawed questions during training to improve the model’s ability to identify and respond appropriately.
  - (b) **Enhanced Feedback Mechanisms:** Utilizing reinforcement learning with human feedback (RLHF) or synthetic feedback to refine the model’s judgment on logically flawed scenarios.

## 1.2. Submitting Final Camera-Ready Copy

The final versions of papers accepted for publication should follow the same format and naming convention as initial submissions, except that author information (names and affiliations) should be given. See Section 2.3.2 for formatting instructions.

The footnote, “Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.” must be modified to “*Proceedings of the 42<sup>nd</sup> International Conference on Machine Learning*, Vancouver, Canada, PMLR 267, 2025. Copyright 2025 by the author(s).”

For those using the **L<sup>A</sup>T<sub>E</sub>X** style file, this change (and others) is handled automatically by simply changing `\usepackage{icml2025}` to

```
\usepackage[accepted]{icml2025}
```

Authors using **Word** must edit the footnote on the first page of the document themselves.

Camera-ready copies should have the title of the paper as running head on each page except the first one. The running title consists of a single line centered above a horizontal rule which is 1 point thick. The running head should be centered, bold and in 9 point type. The rule should be 10 points above the main text. For those using the **L<sup>A</sup>T<sub>E</sub>X** style file, the original title is automatically set as running head using the `fancyhdr` package which is included in the ICML 2025 style file package. In case that the original title exceeds the size restrictions, a shorter form can be supplied by using

```
\icmltitlerunning{...}
```

just before `\begin{document}`. Authors using **Word** must edit the header of the document themselves.

## 2. Format of the Paper

All submissions must follow the specified format.

### 2.1. Dimensions

The text of the paper should be formatted in two columns, with an overall width of 6.75 inches, height of 9.0 inches, and 0.25 inches between the columns. The left margin should be 0.75 inches and the top margin 1.0 inch (2.54 cm). The right and bottom margins will depend on whether you print on US letter or A4 paper, but all final versions must be produced for US letter size. Do not write anything on the margins.

The paper body should be set in 10 point type with a vertical spacing of 11 points. Please use Times typeface throughout the text.

### 2.2. Title

The paper title should be set in 14 point bold type and centered between two horizontal rules that are 1 point thick, with 1.0 inch between the top rule and the top edge of the page. Capitalize the first letter of content words and put the rest of the title in lower case.

---

## 2.3. Author Information for Submission

ICML uses double-blind review, so author information must not appear. If you are using L<sup>A</sup>T<sub>E</sub>X and the `icml2025.sty` file, use `\icmlauthor{...}` to specify authors and `\icmlaffiliation{...}` to specify affiliations. (Read the TeX code used to produce this document for an example usage.) The author information will not be printed unless `accepted` is passed as an argument to the style file. Submissions that include the author information will not be reviewed.

### 2.3.1. SELF-CITATIONS

If you are citing published papers for which you are an author, refer to yourself in the third person. In particular, do not use phrases that reveal your identity (e.g., “in previous work (Langley, 2000), we have shown ...”).

Do not anonymize citations in the reference section. The only exception are manuscripts that are not yet published (e.g., under submission). If you choose to refer to such unpublished manuscripts (Author, 2021), anonymized copies have to be submitted as Supplementary Material via Open-Review. However, keep in mind that an ICML paper should be self contained and should contain sufficient detail for the reviewers to evaluate the work. In particular, reviewers are not required to look at the Supplementary Material when writing their review (they are not required to look at more than the first 8 pages of the submitted document).

### 2.3.2. CAMERA-READY AUTHOR INFORMATION

If a paper is accepted, a final camera-ready copy must be prepared. For camera-ready papers, author information should start 0.3 inches below the bottom rule surrounding the title. The authors’ names should appear in 10 point bold type, in a row, separated by white space, and centered. Author names should not be broken across lines. Unbolded superscripted numbers, starting 1, should be used to refer to affiliations.

Affiliations should be numbered in the order of appearance. A single footnote block of text should be used to list all the affiliations. (Academic affiliations should list Department, University, City, State/Region, Country. Similarly for industrial affiliations.)

Each distinct affiliations should be listed once. If an author has multiple affiliations, multiple superscripts should be placed after the name, separated by thin spaces. If the authors would like to highlight equal contribution by multiple first authors, those authors should have an asterisk placed after their name in superscript, and the term “\*Equal contribution” should be placed in the footnote block ahead of the list of affiliations. A list of corresponding authors and their emails (in the format Full Name <email@domain.com>)

can follow the list of affiliations. Ideally only one or two names should be listed.

A sample file with author names is included in the ICML2025 style file package. Turn on the `[accepted]` option to the stylefile to see the names rendered. All of the guidelines above are implemented by the L<sup>A</sup>T<sub>E</sub>X style file.

## 2.4. Abstract

The paper abstract should begin in the left column, 0.4 inches below the final address. The heading ‘Abstract’ should be centered, bold, and in 11 point type. The abstract body should use 10 point type, with a vertical spacing of 11 points, and should be indented 0.25 inches more than normal on left-hand and right-hand margins. Insert 0.4 inches of blank space after the body. Keep your abstract brief and self-contained, limiting it to one paragraph and roughly 4–6 sentences. Gross violations will require correction at the camera-ready phase.

## 2.5. Partitioning the Text

You should organize your paper into sections and paragraphs to help readers place a structure on the material and understand its contributions.

### 2.5.1. SECTIONS AND SUBSECTIONS

Section headings should be numbered, flush left, and set in 11 pt bold type with the content words capitalized. Leave 0.25 inches of space before the heading and 0.15 inches after the heading.

Similarly, subsection headings should be numbered, flush left, and set in 10 pt bold type with the content words capitalized. Leave 0.2 inches of space before the heading and 0.13 inches afterward.

Finally, subsubsection headings should be numbered, flush left, and set in 10 pt small caps with the content words capitalized. Leave 0.18 inches of space before the heading and 0.1 inches after the heading.

Please use no more than three levels of headings.

### 2.5.2. PARAGRAPHS AND FOOTNOTES

Within each section or subsection, you should further partition the paper into paragraphs. Do not indent the first line of a given paragraph, but insert a blank line between succeeding ones.

You can use footnotes<sup>1</sup> to provide readers with additional information about a topic without interrupting the flow of the paper. Indicate footnotes with a number in the text

---

<sup>1</sup>Footnotes should be complete sentences.

where the point is most relevant. Place the footnote in 9 point type at the bottom of the column in which it appears. Precede the first footnote in a column with a horizontal rule of 0.8 inches.<sup>2</sup>

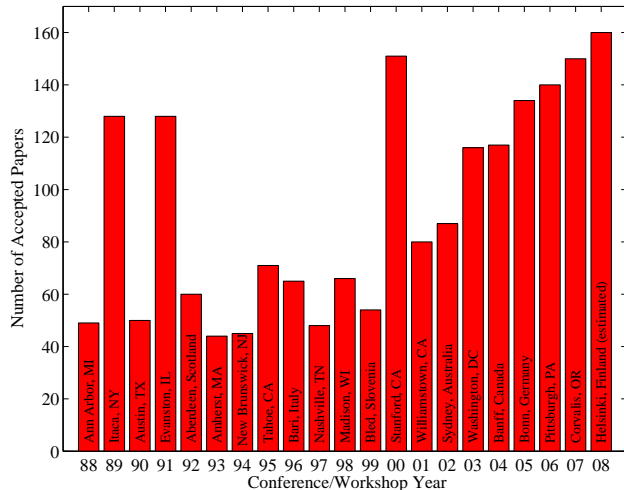


Figure 1. Historical locations and number of accepted papers for International Machine Learning Conferences (ICML 1993 – ICML 2008) and International Workshops on Machine Learning (ML 1988 – ML 1992). At the time this figure was produced, the number of accepted papers for ICML 2008 was unknown and instead estimated.

## 2.6. Figures

You may want to include figures in the paper to illustrate your approach and results. Such artwork should be centered, legible, and separated from the text. Lines should be dark and at least 0.5 points thick for purposes of reproduction, and text should not appear on a gray background.

Label all distinct components of each figure. If the figure takes the form of a graph, then give a name for each axis and include a legend that briefly describes each curve. Do not include a title inside the figure; instead, the caption should serve this function.

Number figures sequentially, placing the figure number and caption *after* the graphics, with at least 0.1 inches of space before the caption and 0.1 inches after it, as in Figure 1. The figure caption should be set in 9 point type and centered unless it runs two or more lines, in which case it should be flush left. You may float figures to the top or bottom of a column, and you may set wide figures across both columns (use the environment `figure*` in  $\text{\LaTeX}$ ). Always place

<sup>2</sup>Multiple footnotes can appear in each column, in the same order as they appear in the text, but spread them across columns and pages if possible.

## Algorithm 1 Bubble Sort

---

**Input:** data  $x_i$ , size  $m$

**repeat**

Initialize  $noChange = true$ .

**for**  $i = 1$  **to**  $m - 1$  **do**

**if**  $x_i > x_{i+1}$  **then**

Swap  $x_i$  and  $x_{i+1}$

$noChange = false$

**end if**

**end for**

**until**  $noChange$  is *true*

---

Table 1. Classification accuracies for naive Bayes and flexible Bayes on various data sets.

DATA SET	NAIVE	FLEXIBLE	BETTER?
BREAST	95.9±0.2	96.7±0.2	✓
CLEVELAND	83.3±0.6	80.0±0.6	×
GLASS2	61.9±1.4	83.8±0.7	✓
CREDIT	74.8±0.5	78.3±0.6	
HORSE	73.3±0.9	69.7±1.0	×
META	67.1±0.6	76.5±0.5	✓
PIMA	75.1±0.6	73.9±0.5	
VEHICLE	44.9±0.6	61.5±0.4	✓

two-column figures at the top or bottom of the page.

## 2.7. Algorithms

If you are using  $\text{\LaTeX}$ , please use the “algorithm” and “algorithmic” environments to format pseudocode. These require the corresponding stylefiles, `algorithm.sty` and `algorithmic.sty`, which are supplied with this package. Algorithm 1 shows an example.

## 2.8. Tables

You may also want to include tables that summarize material. Like figures, these should be centered, legible, and numbered consecutively. However, place the title *above* the table with at least 0.1 inches of space before the title and the same after it, as in Table 1. The table title should be set in 9 point type and centered unless it runs two or more lines, in which case it should be flush left.

Tables contain textual material, whereas figures contain graphical material. Specify the contents of each row and column in the table’s topmost row. Again, you may float tables to a column’s top or bottom, and set wide tables across both columns. Place two-column tables at the top or bottom of the page.

## 2.9. Theorems and such

The preferred way is to number definitions, propositions, lemmas, etc. consecutively, within sections, as shown below.

**Definition 2.1.** A function  $f : X \rightarrow Y$  is injective if for any  $x, y \in X$  different,  $f(x) \neq f(y)$ .

Using Definition 2.1 we immediately get the following result:

**Proposition 2.2.** *If  $f$  is injective mapping a set  $X$  to another set  $Y$ , the cardinality of  $Y$  is at least as large as that of  $X$*

*Proof.* Left as an exercise to the reader.  $\square$

Lemma 2.3 stated next will prove to be useful.

**Lemma 2.3.** *For any  $f : X \rightarrow Y$  and  $g : Y \rightarrow Z$  injective functions,  $f \circ g$  is injective.*

**Theorem 2.4.** *If  $f : X \rightarrow Y$  is bijective, the cardinality of  $X$  and  $Y$  are the same.*

An easy corollary of Theorem 2.4 is the following:

**Corollary 2.5.** *If  $f : X \rightarrow Y$  is bijective, the cardinality of  $X$  is at least as large as that of  $Y$ .*

**Assumption 2.6.** The set  $X$  is finite.

*Remark 2.7.* According to some, it is only the finite case (cf. Assumption 2.6) that is interesting.

## 2.10. Citations and References

Please use APA reference format regardless of your formatter or word processor. If you rely on the  $\text{\LaTeX}$  bibliographic facility, use `natbib.sty` and `icml2025.bst` included in the style-file package to obtain this format.

Citations within the text should include the authors' last names and year. If the authors' names are included in the sentence, place only the year in parentheses, for example when referencing Arthur Samuel's pioneering work (1959). Otherwise place the entire reference in parentheses with the authors and year separated by a comma (Samuel, 1959). List multiple references separated by semicolons (Kearns, 1989; Samuel, 1959; Mitchell, 1980). Use the 'et al.' construct only for citations with three or more authors or after listing all authors to a publication in an earlier reference (Michalski et al., 1983).

Authors should cite their own work in the third person in the initial version of their paper submitted for blind review. Please refer to Section 2.3 for detailed instructions on how to cite your own papers.

Use an unnumbered first-level section heading for the references, and use a hanging indent style, with the first line of

the reference flush against the left margin and subsequent lines indented by 10 points. The references at the end of this document give examples for journal articles (Samuel, 1959), conference publications (Langley, 2000), book chapters (Newell & Rosenbloom, 1981), books (Duda et al., 2000), edited volumes (Michalski et al., 1983), technical reports (Mitchell, 1980), and dissertations (Kearns, 1989).

Alphabetize references by the surnames of the first authors, with single author entries preceding multiple author entries. Order references for the same authors by year of publication, with the earliest first. Make sure that each reference includes all relevant information (e.g., page numbers).

Please put some effort into making references complete, presentable, and consistent, e.g. use the actual current name of authors. If using bibtex, please protect capital letters of names and abbreviations in titles, for example, use `\B}ayesian` or `\L}ipschitz` in your .bib file.

## Accessibility

Authors are kindly asked to make their submissions as accessible as possible for everyone including people with disabilities and sensory or neurological differences. Tips of how to achieve this and what to pay attention to will be provided on the conference website <http://icml.cc/>.

## Software and Data

If a paper is accepted, we strongly encourage the publication of software and data with the camera-ready version of the paper whenever appropriate. This can be done by including a URL in the camera-ready copy. However, **do not** include URLs that reveal your institution or identity in your submission for review. Instead, provide an anonymous URL or upload the material as "Supplementary Material" into the OpenReview reviewing system. Note that reviewers are not required to look at this material when writing their review.

## Acknowledgements

**Do not** include acknowledgements in the initial version of the paper submitted for blind review.

If a paper is accepted, the final camera-ready version can (and usually should) include acknowledgements. Such acknowledgements should be placed at the end of the section, in an unnumbered section that does not count towards the paper page limit. Typically, this will include thanks to reviewers who gave useful comments, to colleagues who contributed to the ideas, and to funding agencies and corporate sponsors that provided financial support.

---

## Impact Statement

Authors are **required** to include a statement of the potential broader impact of their work, including its ethical aspects and future societal consequences. This statement should be in an unnumbered section at the end of the paper (co-located with Acknowledgements – the two may appear in either order, but both must be before References), and does not count toward the paper page limit. In many cases, where the ethical impacts and expected societal implications are those that are well established when advancing the field of Machine Learning, substantial discussion is not required, and a simple statement such as the following will suffice:

“This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.”

The above statement can be used verbatim in such cases, but we encourage authors to think about whether there is content which does warrant further discussion, as this statement will be apparent if the paper is later flagged for ethics review.

## References

- Author, N. N. Suppressed for anonymity, 2021.
- Duda, R. O., Hart, P. E., and Stork, D. G. *Pattern Classification*. John Wiley and Sons, 2nd edition, 2000.
- Kearns, M. J. *Computational Complexity of Machine Learning*. PhD thesis, Department of Computer Science, Harvard University, 1989.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Michalski, R. S., Carbonell, J. G., and Mitchell, T. M. (eds.). *Machine Learning: An Artificial Intelligence Approach, Vol. I*. Tioga, Palo Alto, CA, 1983.
- Mitchell, T. M. The need for biases in learning generalizations. Technical report, Computer Science Department, Rutgers University, New Brunswick, MA, 1980.
- Newell, A. and Rosenbloom, P. S. Mechanisms of skill acquisition and the law of practice. In Anderson, J. R. (ed.), *Cognitive Skills and Their Acquisition*, chapter 1, pp. 1–51. Lawrence Erlbaum Associates, Inc., Hillsdale, NJ, 1981.
- Samuel, A. L. Some studies in machine learning using the game of checkers. *IBM Journal of Research and Development*, 3(3):211–229, 1959.

---

## **A. You *can* have an appendix here.**

You can have as much text here as you want. The main body must be at most 8 pages long. For the final version, one more page can be added. If you want, you can use an appendix like this one.

The `\onecolumn` command above can be kept in place if you prefer a one-column appendix, or can be removed if you prefer a two-column appendix. Apart from this possible change, the style (font size, spacing, margins, page numbering, etc.) should be kept the same as the main body.