Grading Research Report #1

The goal of this assignment is to read a research report that was generated by an AI tool and provide quality ratings across various measures. You will be provided with a link to a report and you should read this report in its entirety

You should read this report through the following lens of perception:

You were provided with an Al assistant tool that you tasked to perform research on the following question: "Does gender role play affect the accuracy on of language models on answering math questions?". Given the question that you provided, the Al assistant tool performed its own literature search, experimentation, performed its own coding, executed the code, collected data, conducted an analysis, and wrote the presented research report. The goal of this assistant is not to perform research for you (automate you task) but instead to provide a foundation for you to accelerate your own research. You should be asking: is what this Al assistant produced useful for me to build off of instead comparing it to what a human would perform.

Link to the research

report: https://drive.google.com/file/d/19vjnzgbsrkiHL50lxbCHZ_uU20jpzlxx/view? usp=sharing

Once you have read the paper please answer the question below:

* Indicates required question

1. Let's assume you were provided with an AI assistant tool that you tasked to perform research on the following question: "Does gender role play affect the accuracy on of language models on answering math questions?"

What is your perception of the quality of the <u>experimental results</u> presented in this report?

Please provide a rating 1-5, with the following rating descriptions:

- 1 Very Low Quality
- 2 Low Quality
- 3 Medium Quality
- 4 High Quality
- 5 Very High Quality



2. Let's assume you were provided with a research paper answering the following question: "Does gender role play affect the accuracy on of language models on answering math questions?"

What is your perception of the quality of the <u>research report writing quality</u> presented in this report?

Please provide a rating 1-5, with the following rating descriptions:

- 1 Very Low Quality
- 2 Low Quality
- 3 Medium Quality
- 4 High Quality
- 5 Very High Quality



What is your perception of the <u>usefulness of the Al assistant tool</u> presented in this report?

Please provide a rating 1-5, with the following rating descriptions:

- 1 Very Low Usefulness
- 2 Low Usefulness
- 3 Medium Usefulness
- 4 High Usefulness
- 5 Very High Usefulness



Review

Now assume you are a reviewer at NeurIPS 2025 and are reviewing a machine learning paper.

Please provide the following ratings from this perspective.

*

Quality: Is the submission technically sound? Are claims well supported (e.g., by theoretical analysis or experimental results)? Are the methods used appropriate? Is this a complete piece of work or work in progress? Are the authors careful and honest about evaluating both the strengths and weaknesses of their work

- 1 Low Quality
- 2 Medium Quality
- 3 High Quality
- 4 Very High Quality



5. Let's assume you were provided with a research paper answering the following question: "Does gender role play affect the accuracy on of language models on answering math questions?"

Clarity: Is the submission clearly written? Is it well organized? (If not, please make constructive suggestions for improving its clarity.) Does it adequately inform the reader? (Note that a superbly written paper provides enough information for an expert reader to reproduce its results.)

Please provide a rating 1-4, with the following rating descriptions:

- 1 Low Clarity
- 2 Medium Clarity
- 3 High Clarity
- 4 Very High Clarity



Significance: Are the results important? Are others (researchers or practitioners) likely to use the ideas or build on them? Does the submission address a difficult task in a better way than previous work? Does it advance the state of the art in a demonstrable way? Does it provide unique data, unique conclusions about existing data, or a unique theoretical or experimental approach?

Please provide a rating 1-4, with the following rating descriptions:

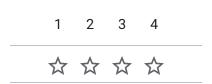
- 1 Low Significance
- 2 Medium Significance
- 3 High Significance
- 4 Very High Significance



7. Let's assume you were provided with a research paper answering the following question: "Does gender role play affect the accuracy on of language models on answering math questions?"

Soundness: Please assign the paper a numerical rating on the following scale to indicate the soundness of the technical claims, experimental and research methodology and on whether the central claims of the paper are adequately supported with evidence.

- 4: excellent
- 3: good
- 2: fair
- 1: poor



Presentation: Please assign the paper a numerical rating on the following scale to indicate the quality of the presentation. This should take into account the writing style and clarity, as well as contextualization relative to prior work.

- 4: excellent
- 3: good
- 2: fair
- 1: poor



9. Let's assume you were provided with a research paper answering the following question: "Does gender role play affect the accuracy on of language models on answering math questions?"

Contribution: Please assign the paper a numerical rating on the following scale to indicate the quality of the overall contribution this paper makes to the research area being studied. Are the questions being asked important? Does the paper bring a significant originality of ideas and/or execution? Are the results valuable to share with the broader NeurIPS community.

- 4: excellent
- 3: good
- 2: fair
- 1: poor



Overall: Please provide an "overall score" for this submission. Choices:

- 10: Award quality
- 9: Very Strong Accept
- 8: Strong Accept
- 7: Accept
- 6: Weak Accept
- 5: Borderline accept
- 4: Borderline reject
- 3: Reject
- 2: Strong Reject
- 1: Very Strong Reject

1 2 3 4 5 6 7 8 9 10

Confidence: Please provide a "confidence score" for your assessment of this submission to indicate how confident you are in your evaluation.

Choices:

- 5: You are absolutely certain about your assessment. You are very familiar with the related work and checked the math/other details carefully.
- 4: You are confident in your assessment, but not absolutely certain. It is unlikely, but not impossible, that you did not understand some parts of the submission or that you are unfamiliar with some pieces of related work.
- 3: You are fairly confident in your assessment. It is possible that you did not understand some parts of the submission or that you are unfamiliar with some pieces of related work. Math/other details were not carefully checked.
- 2: You are willing to defend your assessment, but it is quite likely that you did not understand the central parts of the submission or that you are unfamiliar with some pieces of related work. Math/other details were not carefully checked.
- 1: Your assessment is an educated guess. The submission is not in your area or the submission was difficult to understand. Math/other details were not carefully checked.

1	2	3	4	5
☆	☆	☆	☆	☆

12. Let's assume you were provided with a research paper answering the following question: "Does gender role play affect the accuracy on of language models on answering math questions?"

"Decision": A decision that has to be one of the following: Accept, Reject.

Mark only one oval.

