A Systematic Review of Algorithmic Red Teaming Methodologies for AI Assurance and Security

Shaurya Jauhari^{1*}, Shruti Srivastava^{1†} and Kiranmayee Janardhan^{1†}

^{1*}Responsible AI Office, Infosys Limited, Electronic City, Bangalore, 560100, Karnataka, India.

*Corresponding author(s). E-mail(s): shaurya.jauhari@infosys.com;
Contributing authors: shruti.srivastava03@infosys.com;
kiranmayee.j@infosys.com;

[†]These authors contributed equally to this work.

11 Abstract

While Generative AI has unlocked creative potential across numerous disciplines and advanced capabilities like assisted data generation, its rapid adoption brings significant security concerns to the forefront. These issues are as critical as those faced by earlier predictive AI models, with a growing body of literature demonstrating that Large Language Models (LLMs) can be manipulated into generating harmful or inappropriate content. To address this, red teaming has become a standard practice to systematically identify such vulnerabilities and assess a model's resilience. This review provides a systematic documentation of prominent algorithmic red teaming approaches and outlines the desirable characteristics that effective LLM red teaming frameworks should exhibit.

Keywords: Algorithmic Detection, Automation, LLM Security, PRISMA, Red Teaming, Systematic Review

4 1 Introduction

10

12

14

15

16

17

18

19

20

21

22

23

Artificial Intelligence (AI), a term coined by John McCarthy in 1956 [1–5], is a domain having deeper roots in history that date back to mid 20th century, with Warren McCulloch and Walter Pitts developing the first mathematical model of a neural network in 1943 [6, 7]. In fact, John would later invent List Programming Language (LISP)

to corroborate AI with efficient incorporation of symbolic information. Later, in 1950,
Alan Turing devised the revered <u>Turing Test</u> to benchmark machine intelligence [8].
While the mathematical pedigree of the machine learning algorithms is centuries old,
literally, it was only with the emergence of huge volumes of data that the shelved algorithms started making sense. The models when trained with necessary and sufficient quantum of data, are rendered potent to manifest "intellect".

Models can be both predictive and now, generative. While the former have the capacity to parse through the input data and label it based on the training data patterns, the latter go a step further to create novel data instances. This has proven revolutionary, especially with Generative Pre-trained Transformer (GPT) models breaking the ice [9], although Google was first to arrive at the scene with transformer models [10]. Nonetheless, while the usefulness of these, and many others that have flooded the market, have seldom been debated, what remains a challenge still is keeping the data and the concerned applications secure. Security is undeniably indispensable.

Red Teaming has been a staple feature in the cyber-security domain, as a strategy to safeguard against any potential espionage that adversaries could attempt [11]. In a controlled, sand-boxed environment a series of manual and/or automated attacks are simulated on the target application to elicit signs of fragility [12]. Actually, red teaming in military and intelligence operations, originally, found way to penetrate enterprise settings and that is how it flourished in civil space [13].

9 2 Results

35

37

41

44

Sample body text. Sample body text.

₅₂ 3 This is an example for first level head—section head

3.1 This is an example for second level head—subsection head

3.1.1 This is an example for third level head—subsubsection head

Sample body text. Sample body text.

$_{57}$ 4 Equations

Equations in IATEX can either be inline or on-a-line by itself ("display equations"). For inline equations use the \$...\$ commands. E.g.: The equation $H\psi = E\psi$ is written via the command \$H \psi = E \psi\$.

For display equations (with auto generated equation numbers) one can use the equation or align environments:

$$\|\tilde{X}(k)\|^{2} \leq \frac{\sum_{i=1}^{p} \|\tilde{Y}_{i}(k)\|^{2} + \sum_{j=1}^{q} \|\tilde{Z}_{j}(k)\|^{2}}{p+q}.$$
 (1)

63 where,

$$D_{\mu} = \partial_{\mu} - ig \frac{\lambda^{a}}{2} A^{a}_{\mu}$$

$$F^{a}_{\mu\nu} = \partial_{\mu} A^{a}_{\nu} - \partial_{\nu} A^{a}_{\mu} + g f^{abc} A^{b}_{\mu} A^{a}_{\nu}$$

$$(2)$$

Notice the use of \nonumber in the align environment at the end of each line, except the last, so as not to produce equation numbers on lines where no equation numbers are required. The \label{} command should only be used at the last line of an align environment where \nonumber is not used.

$$Y_{\infty} = \left(\frac{m}{\text{GeV}}\right)^{-3} \left[1 + \frac{3\ln(m/\text{GeV})}{15} + \frac{\ln(c_2/5)}{15}\right]$$
 (3)

The class file also supports the use of \mathcal{R} , \mathscr{} and \mathcal{} commands. As such \mathbb{R}, \mathscr{R} and \mathcal{R} produces \mathbb{R} , \mathcal{R} and \mathcal{R} respectively (refer Subsubsection 3.1.1).

$_{71}$ 5 Tables

Tables can be inserted via the normal table and tabular environment. To put footnotes inside tables you should use \footnotetext[]{...} tag. The footnote appears just below the table itself (refer Tables 1 and 2). For the corresponding footnotemark use \footnotemark[...]

Table 1 Caption text

| Column 1 | Column 2 | Column 3 | Column 4 |
|----------|----------|--|---------------------|
| row 1 | data 1 | $\begin{array}{c} \text{data 2} \\ \text{data 5}^1 \\ \text{data 8} \end{array}$ | data 3 |
| row 2 | data 4 | | data 6 |
| row 3 | data 7 | | data 9 ² |

Source: This is an example of table footnote. This is an example of table footnote.

- The input format for the above table is as follows:
- 77 \begin{table}[<placement-specifier>]
- 78 \caption{<table-caption>}\label{<table-label>}%
- 79 \begin{tabular}{@{}1111@{}}

 $^{^1{\}rm Example}$ for a first table footnote. This is an example of table footnote.

 $^{^2\}mathrm{Example}$ for a second table footnote. This is an example of table footnote.

```
80 \toprule
```

- 81 Column 1 & Column 2 & Column 3 & Column 4\\
- 82 \midrule
- 3 row 1 & data 1 & data 2 & data 3 \\
- row 2 & data 4 & data $5\footnotemark[1]$ & data $6\$
- s row 3 & data 7 & data 8 & data 9\footnotemark[2]\\
- 86 \botrule
- 87 \end{tabular}
- \footnotetext{Source: This is an example of table footnote.
- 89 This is an example of table footnote.}
- o \footnotetext[1]{Example for a first table footnote.
- 91 This is an example of table footnote.}
- 92 \footnotetext[2]{Example for a second table footnote.
- 93 This is an example of table footnote.}
- 94 \end{table}

Table 2 Example of a lengthy table which is set to full textwidth

| | | Element 1 | 1 | | Element 2 | 2^2 |
|------------------------|----------------|-----------------|-------------------------------|----------------|-----------------|---------------------------------|
| Project | Energy | σ_{calc} | σ_{expt} | Energy | σ_{calc} | σ_{expt} |
| Element 3 Element 4 | 990 A 500 A | 1168 961 | 1547 ± 12 922 ± 10 | 780 A 900 A | 1166 1268 | 1239 ± 100 1092 ± 40 |

Note: This is an example of table footnote. This is an example of table footnote this is an example of table footnote this is an example of table footnote.

In case of double column layout, tables which do not fit in single column width should be set to full text width. For this, you need to use \begin{table*} ... \end{table*} instead of \begin{table} ... \end{table} environment. Lengthy tables which do not fit in textwidth should be set as rotated table. For this, you need to use \begin{sidewaystable} ... \end{sidewaystable} instead of \begin{table*} ... \end{table*} environment. This environment puts tables rotated to single column width. For tables rotated to double column width, use \begin{sidewaystable*} ... \end{sidewaystable*}... \end{sidewaystable*}.

6 Figures

103

As per the LaTeX standards you need to use eps images for LaTeX compilation and pdf/jpg/png images for PDFLaTeX compilation. This is one of the major difference between LaTeX and PDFLaTeX. Each image should be from a single input .eps/vector image file. Avoid using subfigures. The command for inserting images for LaTeX and

 $^{^1}$ Example for a first table footnote.

²Example for a second table footnote.

Table 3 Tables which are too long to fit, should be written using the "sidewaystable" environment as shown here

| | | Element 1 ¹ | | | $ m Element^2$ | |
|------------|--------|------------------------|-----------------|--------|-----------------|-----------------|
| Projectile | Energy | σ_{calc} | σ_{expt} | Energy | σ_{calc} | σ_{expt} |
| Element 3 | 990 A | 1168 | 1547 ± 12 | 780 A | 1166 | 1239 ± 100 |
| Element 4 | 500 A | 961 | 922 ± 10 | 900 A | 1268 | 1092 ± 40 |
| Element 5 | 990 A | 1168 | 1547 ± 12 | 780 A | 1166 | 1239 ± 100 |
| Element 6 | 500 A | 961 | 922 ± 10 | 900 A | 1268 | 1092 ± 40 |
| | | | | | | |

Note: This is an example of table footnote this is an example of table footnote this is an example of table footnote this is an example of table footnote.

¹This is an example of table footnote.

```
PDFLaTeX can be generalized. The package used to insert images in LaTeX/PDFLaTeX is the graphicx package. Figures can be inserted via the normal figure environment as shown in the below example:
```

```
\text{\text{begin{figure} [<placement-specifier>]} \text{\text{centering} \text{\text{includegraphics{<eps-file>}} \text{\text{caption{<figure-caption>}\label{<figure-label>}} \text{\text{end{figure}}}
```

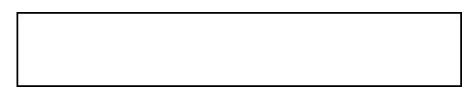


Fig. 1 This is a widefig. This is an example of long caption this is an example of long caption this is an example of long caption

In case of double column layout, the above format puts figure captions/images to single column width. To get spanned images, we need to provide \begin{figure*}
... \end{figure*}.

For sample purpose, we have included the width of images in the optional argument of \includegraphics tag. Please ignore this.

7 Algorithms, Program codes and Listings

Packages algorithm, algorithmicx and algorithmed are used for setting algorithms in LATEX using the format:

```
124 \begin{algorithm}
125 \caption{<alg-caption>}\label{<alg-label>}
126 \begin{algorithmic}[1]
127 . . .
128 \end{algorithmic}
129 \end{algorithm}
```

You may refer above listed package documentations for more details before setting algorithm environment. For program codes, the "verbatim" package is required and the command to be used is \begin{verbatim} . . . \end{verbatim}.

```
Similarly, for listings, use the listings package. \begin{lstlisting} ...
133
    \end{lstlisting} is used to set environments similar to verbatim environment.
134
    Refer to the lstlisting package documentation for more details.
135
       A fast exponentiation procedure:
      for i := 1 to 10 step 1 do
            expt(2,i);
            newline() od
                                                Comments will be set flush to the right margin
140
    where
141
    proc expt(x,n) \equiv
142
      z := 1;
143
      do if n=0 then exit fi;
144
          do if odd(n) then exit fi;
145
              comment: This is a comment statement;
146
147
              n := n/2; \quad x := x * x \text{ od};
          \{ n > 0 \};
          n:=n-1\,;\ z:=z*x\ \operatorname{od}\,;
149
       print(z).
150
    end
151
```

Algorithm 1 Calculate $y = x^n$

```
Require: n \ge 0 \lor x \ne 0
Ensure: y = x^n
 1: y \Leftarrow 1
 2: if n < 0 then
         X \Leftarrow 1/x
 3:
          N \Leftarrow -n
 4:
 5: else
          X \Leftarrow x
 6:
         N \Leftarrow n
 7:
 8: end if
     while N \neq 0 do
          if N is even then
10:
              X \Leftarrow X \times X
11:
              N \Leftarrow N/2
12:
          else[N \text{ is odd}]
13:
              y \Leftarrow y \times X
14:
              N \Leftarrow N - 1
15:
          end if
16:
17: end while
```

```
for i:=maxint to 0 do
    begin
    { do nothing }
    end;
    Write('Case-insensitive-');
    Write('Pascal-keywords.');
```

8 Cross referencing

Environments such as figure, table, equation and align can have a label declared via the \label{#label} command. For figures and table environments use the \label{} command inside or just below the \caption{} command. You can then use the \ref{#label} command to cross-reference them. As an example, consider the label declared for Figure 1 which is \label{fig1}. To cross-reference it, use the command Figure \ref{fig1}, for which it comes up as "Figure 1".

To reference line numbers in an algorithm, consider the label declared for the line number 2 of Algorithm 1 is \label{algln2}. To cross-reference it, use the command \ref{algln2} for which it comes up as line 2 of Algorithm 1.

8.1 Details on reference citations

Standard IATEX permits only numerical citations. To support both numerical and author-year citations this template uses natbib IATEX package. For style guidance please refer to the template user manual.

Here is an example for \cite{...}: [?]. Another example for \citep{...}: [?]. For author-year citation mode, \cite{...} prints Jones et al. (1990) and \citep{...} prints (Jones et al., 1990).

9 Examples for theorem like environments

For theorem like environments, we require amsthm package. There are three types of predefined theorem styles exists—thmstyleone, thmstyletwo and thmstylethree

| thmstyleone | Numbered, theorem head in bold font and theorem |
|---------------|--|
| | text in italic style |
| thmstyletwo | Numbered, theorem head in roman font and theorem |
| | text in italic style |
| thmstylethree | Numbered, theorem head in bold font and theorem |
| | text in roman style |

For mathematics journals, theorem styles can be included as shown in the following examples:

Theorem 1 (Theorem subhead) Example theorem text. Example theorem text.

Sample body text. Sample body text. Sample body text. Sample body text. Sample body text. Sample body text. Sample body text.

Proposition 2 Example proposition text. Example proposition text.

Sample body text. Sample body text. Sample body text. Sample body text. Sample body text. Sample body text. Sample body text.

Example 1 Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem.

Sample body text. Sample body text.

Remark 1 Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis,
 molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at,
 accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat
 lorem.

Sample body text. Sample body text.

Definition 1 (Definition sub head) Example definition text. Example definition text.

Additionally a predefined "proof" environment is available: \begin{proof} ... \end{proof}. This prints a "Proof" head in italic font style and the "body text" in roman font style with an open square at the end of each proof environment.

Proof Example for proof text. Example for proof text.

Sample body text. Sample body text. Sample body text. Sample body text. Sample body text. Sample body text. Sample body text.

Proof of Theorem 1 Example for proof text. D

For a quote environment, use \begin{quote}...\end{quote}

Quoted text example. Aliquam porttitor quam a lacus. Praesent vel arcu ut tortor cursus volutpat. In vitae pede quis diam bibendum placerat. Fusce elementum convallis neque. Sed dolor orci, scelerisque ac, dapibus nec, ultricies ut, mi. Duis nec dui quis leo sagittis commodo.

Sample body text. Sample body text (refer Table 3).

4 10 Methods

Topical subheadings are allowed. Authors must ensure that their Methods section includes adequate experimental and characterization data necessary for others in the field to reproduce their work. Authors are encouraged to include RIIDs where appropriate.

Ethical approval declarations (only required where applicable) Any article reporting experiment/s carried out on (i) live vertebrate (or higher invertebrates), (ii) humans or (iii) human samples must include an unambiguous statement within the methods section that meets the following requirements:

- 233 1. Approval: a statement which confirms that all experimental protocols were approved by a named institutional and/or licensing committee. Please identify the approving body in the methods section
 - 2. Accordance: a statement explicitly saying that the methods were carried out in accordance with the relevant guidelines and regulations
 - 3. Informed consent (for experiments involving humans or human tissue samples): include a statement confirming that informed consent was obtained from all participants and/or their legal guardian/s

If your manuscript includes potentially identifying patient/participant information, or if it describes human transplantation research, or if it reports results of a clinical trial then additional information will be required. Please visit (https://www.nature.com/nature-research/editorial-policies) for Nature Portfolio journals, (https://www.springer.com/gp/authors-editors/journal-author/journal-author-helpdesk/publishing-ethics/14214) for Springer Nature journals, or (https://www.biomedcentral.com/getpublished/editorial-policies#ethics+and+consent) for BMC.

₉ 11 Discussion

Discussions should be brief and focused. In some disciplines use of Discussion or 'Conclusion' is interchangeable. It is not mandatory to use both. Some journals prefer a section 'Results and Discussion' followed by a section 'Conclusion'. Please refer to Journal-level guidance for any specific requirements.

²⁵⁴ 12 Conclusion

Conclusions may be used to restate your hypothesis or research question, restate your
 major findings, explain the relevance and the added value of your work, highlight any
 limitations of your study, describe future directions for research and recommendations.
 In some disciplines use of Discussion or 'Conclusion' is interchangeable. It is

258 In some disciplines use of Discussion of Conclusion is interchangeable. It is 259 not mandatory to use both. Please refer to Journal-level guidance for any specific 260 requirements.

Supplementary information. If your article has accompanying supplementary file/s please state so here.

Authors reporting data from electrophoretic gels and blots should supply the full unprocessed scans for key as part of their Supplementary information. This may be requested by the editorial team/s if it is missing.

Please refer to Journal-level guidance for any specific requirements.

Acknowledgements. Acknowledgements are not compulsory. Where included they should be brief. Grant or contribution numbers may be acknowledged.

Please refer to Journal-level guidance for any specific requirements.

70 Declarations

Some journals require declarations to be submitted in a standardised format. Please check the Instructions for Authors of the journal to which you are submitting to see if you need to complete this section. If yes, your manuscript must contain the following sections under the heading 'Declarations':

Funding

275

279

281

- Conflict of interest/Competing interests (check journal-specific guidelines for which heading to use)
- \bullet Ethics approval and consent to participate
 - Consent for publication
 - Data availability
 - Materials availability
- Code availability
 - Author contribution

If any of the sections are not relevant to your manuscript, please include the heading and write 'Not applicable' for that section.

- 286 Editorial Policies for:
- Springer journals and proceedings: https://www.springer.com/gp/editorial-policies
- Nature Portfolio journals: https://www.nature.com/nature-research/editorial-policies
- Scientific Reports: https://www.nature.com/srep/journal-policies/editorial-policies
- 290 BMC journals: https://www.biomedcentral.com/getpublished/editorial-policies

291 Appendix A Section title of first appendix

An appendix contains supplementary information that is not an essential part of the text itself but which may be helpful in providing a more comprehensive understanding of the research problem or it is information that is too cumbersome to be included in the body of the paper.

References

- ²⁹⁷ [1] McCarthy, J., Minsky, M.L., Rochester, N., Shannon, C.E.: A proposal for the dartmouth summer research project on artificial intelligence, august 31, 1955. AI magazine **27**(4), 12–12 (2006)
- [2] McCarthy, J., Hayes, P.J.: Some philosophical problems from the standpoint of artificial intelligence. In: Readings in Artificial Intelligence, pp. 431–450. Elsevier, ??? (1981)
- ³⁰³ [3] McCarthy, J., et al.: What is artificial intelligence (2007)
- McCarthy, J.: Generality in artificial intelligence. Communications of the ACM **30**(12), 1030–1035 (1987)
- McCarthy, J.: Epistemological problems of artificial intelligence. In: Readings in Artificial Intelligence, pp. 459–465. Elsevier, ??? (1981)
- McCulloch, W.S., Pitts, W.: A logical calculus of the ideas immanent in nervous activity. The bulletin of mathematical biophysics **5**(4), 115–133 (1943)
- Pitts, W., McCulloch, W.S.: How we know universals the perception of auditory and visual forms. The Bulletin of mathematical biophysics **9**(3), 127–147 (1947)
- [8] Moor, J.H.: An analysis of the turing test. Philosophical Studies: An International Journal for Philosophy in the Analytic Tradition **30**(4), 249–257 (1976)
- [9] Floridi, L., Chiriatti, M.: Gpt-3: Its nature, scope, limits, and consequences. Minds
 and machines 30(4), 681–694 (2020)

- I0] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N.,
 Kaiser, L., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems 30 (2017)
- [11] Oakley, J.G.: Professional Red Teaming: Conducting Successful Cybersecurity
 Engagements. Apress, ??? (2019)
- ³²¹ [12] Von Solms, R., Van Niekerk, J.: From information security to cyber security. computers & security **38**, 97–102 (2013)
- 323 [13] Jauhari, S.: Algorithmic red teaming for llms: Why should we care? Available at 324 SSRN 5521758