

## The TRIPOD-LLM Statement: A Targeted Guideline For Reporting Large Language Models Use

**Supplementary Table 3:** TRIPOD-LLM Expanded Checklist (Explanation and Elaboration Light)

Section	Item	Checklist Item	Research Design	LLM Task
Title				
Title	1	<p><b>Identify the study as developing, fine-tuning, and/or evaluating the performance of an LLM, specifying the task, the target population, and the outcome to be predicted.</b></p> <ul style="list-style-type: none"> <li>- Informative titles help with the identification of LLM-based studies by potential readers and also systematic reviewers</li> <li>- Report an informative title that provides important information about the target population and the outcome</li> </ul>	All	All
Abstract				
Abstract	2	<p><b>See TRIPOD-LLM Abstract</b></p> <ul style="list-style-type: none"> <li>- Report an abstract addressing each item in the TRIPOD-LLM for Abstracts checklist</li> </ul>	All	All
Introduction				
Background	3a	<p><b>Explain the healthcare context / use case (e.g., administrative, diagnostic, therapeutic, clinical workflow) and rationale for developing or evaluating the LLM, including references to existing approaches and models.</b></p> <ul style="list-style-type: none"> <li>- Describe the healthcare setting or use case where the LLM is intended to be used.</li> <li>- Where existing approaches or LLMs are available, provide a clear justification for developing a new LLM.</li> <li>- For studies evaluating an existing model, provide the rationale for the evaluation and references to all models being evaluated.</li> <li>- For <i>de novo</i> LLM development and LLM methods studies, the precise healthcare context/use cases may not be determined. In this case, provide examples of potential future healthcare contexts/use cases.</li> </ul>	All	All
	3b	<p><b>Describe the target population and the intended use of the LLM in the context of the care pathway, including its intended users in current gold standard practices (e.g., healthcare professionals, patients, public, or administrators).</b></p> <ul style="list-style-type: none"> <li>- Describe who the target population is for the developed or evaluated LLM, such as people of a certain age, in a specific country, or with a specific disease.</li> </ul>	E H	All

		<ul style="list-style-type: none"> <li>- Describe the intended purpose of the LLM, including the clinical decision or guidance the LLM is intended to support (e.g., referral for further testing or hospital admission, triage, starting a treatment, patient portal messaging, billing) and the point in the care pathway where the LLM is intended to be used.</li> <li>- Describe who the intended users of the LLM are, and whether the LLM is for healthcare professionals, patients, public, or other stakeholders.</li> <li>- Explain the current gold standard practices that this LLM is seeking to interact with or replace.</li> </ul>		
Objectives	4	<p><b>Specify the study objectives, including whether the study describes the initial development, fine-tuning, or validation of an LLM (or multiple stages).</b></p> <ul style="list-style-type: none"> <li>- Provide an explicit statement of all objectives of the study, describing whether the study is developing an LLM, fine-tuning or otherwise adjusting an existing LLM, incorporating an existing LLM within a new informatics pipeline or framework, evaluating the performance of an LLM, or covering multiple stages.</li> </ul>	All	All

**Methods**

Data	5a	<p><b>Describe the sources of data separately for the training, tuning, and/or evaluation datasets and the rationale for using these data (e.g., web corpora, clinical research/trial data, EHR data).</b></p> <ul style="list-style-type: none"> <li>- Provide transparency about the data sources used, including whether the data are, for example, from specific web sources, a randomized trial, a registry or from electronic routine healthcare records</li> <li>- Specify whether the study is using existing data or is prospectively collecting new data for the purpose of LLM updating, finetuning or evaluation</li> <li>- Where existing data are being used (i.e., they were originally collected for a different purpose), provide the rationale for using these data, and comment on the suitability (particularly if data are being used from a different setting, country, and/or clinical population to the intended target population) and representativeness of these data with respect to the intended target population and context</li> <li>- If any synthetic data have been used, then provide reasons as to why, and provide all details on how the synthetic data have been created (and code, see item 14f) and used in the study</li> </ul>	All	All
	5b	<p><b>Describe the relevant data points and provide a quantitative and qualitative description of their distribution and other relevant descriptors of the dataset (e.g., source, languages, countries of origin)</b></p> <ul style="list-style-type: none"> <li>- Offer a comprehensive understanding of the dataset used, relevant metadata, languages, and breakdown of characteristics.</li> <li>- Include both quantitative and qualitative descriptions of the data.</li> </ul>	All	All

	5c	<p><b>Specifically state the date of the oldest and newest item of text used in the development process (training, fine-tuning, reward modeling) and in the evaluation datasets.</b></p> <ul style="list-style-type: none"> <li>- Ensure the temporal relevance and validity of the data used for training and/or evaluation.</li> <li>- Provide dates for the text items used in different stages of development and evaluation.</li> <li>- For studies using existing LLMs, provide reference(s) to this information if provided by the original developers or state that this information is not available.</li> </ul>	All	All
	5d	<p><b>Describe any data pre-processing and quality checking, including whether this was similar across text corpora, institutions, and relevant socio-demographic groups.</b></p> <ul style="list-style-type: none"> <li>- If data cleaning is performed e.g. from raw EHR notes, describe any data cleaning steps. This includes transformations of raw data, data quality checks, or translation. All code used for data cleaning should be made available (see item 14e).</li> <li>- Report any efforts in mitigating biased or false content in training.</li> <li>- Report feature selection techniques, if any.</li> <li>- If the data pre-processing/data cleaning steps are extensive, consider reporting this information in the supplementary material.</li> </ul>	All	All
	5e	<p><b>Describe how missing and imbalanced data were handled and provide reasons for omitting any data.</b></p> <ul style="list-style-type: none"> <li>- If the data used are linked with other data or have the potential for missingness (e.g., when extracted from EHR), report any missingness overall and across groups.</li> <li>- If individuals' data have been omitted due to missing values, this should be reported, and reasons given. Note that this is generally not applicable for LLM pretraining.</li> </ul>	All	All
Analytical Methods	6a	<p><b>Report the LLM name, version, and last date of training.</b></p> <ul style="list-style-type: none"> <li>- Given the rapid pace of the field, clear details about the type and version of model used aid in fair comparison across different studies.</li> <li>- For studies using existing LLMs, provide reference(s) to this information if provided by the original developers or state that this information is not available.</li> </ul>	All	All
	6b	<p><b>Report details of LLM development process, such as LLM architecture, training, fine-tuning procedures, and alignment strategy (e.g., reinforcement learning, direct preference optimization, etc.) and alignment goals (e.g., helpfulness, honesty, harmlessness, etc.).</b></p> <ul style="list-style-type: none"> <li>- Outline the complete development process and alignment strategies that were implemented in this study, or point to a study that describes this process.</li> <li>- For any fine-tuning approach, provide details on hyperparameter search and settings, and type of fine-tuning (e.g. full fine-tuning, parameter-efficient fine-tuning strategies).</li> <li>- Specify alignment goals for the LLM and what instructions were given to any labelers involved with the alignment process.</li> </ul>	M D	All

	6c	<p><b>Report details of how text was generated using the LLM, including any prompt engineering and inference settings (e.g., seed, temperature, max token length, penalties), as relevant.</b></p> <ul style="list-style-type: none"> <li>- Describe the model architecture and configuration.</li> <li>- Include details on inference settings such as parameters that control generation, including how these settings were arrived at (e.g., type of sampling used, beam-search).</li> <li>- Provide details on any use of constrained decoding, and any post-processing applied to generated text.</li> </ul>	M D E	All
	6d	<p><b>Specify the initial and post-processed output of the LLM (e.g., probabilities, classification, unstructured text).</b></p> <ul style="list-style-type: none"> <li>- Specify whether the outputs are probabilities, classifications, or unstructured text.</li> <li>- Explain how the initial outputs are transformed or refined in the post-processing stage. All code used for post-processing should be made available (see item 14e).</li> <li>- If the post-processing steps are extensive, consider reporting this information in the supplementary material.</li> </ul>	All	All
	6e	<p><b>Provide details and rationale for any classification and, if applicable, how the probabilities were determined and thresholds identified.</b></p> <ul style="list-style-type: none"> <li>- Describe the process and criteria for categorizing outputs into different classes or groups.</li> <li>- Specify the algorithms or formulas used to derive probability estimates.</li> <li>- Provide a rationale for the chosen thresholds, referencing literature, clinical guidelines, statistical considerations, or ad-hoc decisions.</li> </ul>	All	C OF
LLM Output	7a	<p><b>Include metrics that capture the quality of generative outputs, such as consistency, relevance, similarity, and accuracy, compared to gold standards.</b></p> <ul style="list-style-type: none"> <li>- Given the stochastic nature of LLMs, metrics like consistency, relevance, similarity, and accuracy aid in providing improved characterisation of the results.</li> <li>- Explain how the generative outputs are measured against established benchmarks or reference standards.</li> <li>- Define what gold standard was used or what algorithms or scores derived such metrics.</li> <li>- Provide details of how consistency is measured, e.g. variability to different prompt variations.</li> </ul>	All	QA IR DG SS MT
	7b	<p><b>Report the outcome metrics' relevance to downstream task at deployment time and correlation of metric to human evaluation of the text for the intended use.</b></p> <ul style="list-style-type: none"> <li>- Describe how the outcome metrics are relevant to the real-world application of the LLM.</li> <li>- If human evaluation is carried out, explain how these metrics correlate with human assessments of the text, ensuring the outputs meet user expectations and requirements.</li> </ul>	E H	All
	7c	<p><b>Clearly define the outcome, how the LLM predictions were calculated (e.g., formula, code, object, API), the date of inference for closed-source LLMs, and evaluation metrics.</b></p>	E H	All

		<ul style="list-style-type: none"> <li>- Describe the methodology used for generating LLM output. Include details such as whether a specific algorithm, codebase, software object, or API was used.</li> <li>- Closed-source LLMs may be updated without changes in the named versioning. To enable fair comparisons, report the date of inference for closed-source LLMs.</li> </ul>		
	7d	<p><b>If outcome assessment requires subjective interpretation, describe the qualifications of the assessors, any instructions provided, relevant information on demographics of the assessors, and inter-assessor agreement.</b></p> <ul style="list-style-type: none"> <li>- Provide information on the assessors’ professional background and expertise relevant to the task.</li> <li>- Describe the guidelines and criteria provided to the assessors for the evaluation process.</li> <li>- Include information about the assessors’ demographics to ensure diversity and representativeness.</li> <li>- Report the level of agreement among the assessors using appropriate statistical measures.</li> </ul>	All	All
	7e	<p><b>Specify how performance was compared to other LLMs, humans, and other benchmarks or standards.</b></p> <ul style="list-style-type: none"> <li>- Explain the process and criteria for comparing the LLM’s performance with other models and how these are and are not fair comparisons.</li> <li>- Detail how LLM performance was compared to humans and any differences in the generation and evaluation process between the two groups</li> </ul>	All	All
Annotation	8a	<p><b>If annotation was done, report how text was labeled, including providing specific annotation guidelines with examples.</b></p> <ul style="list-style-type: none"> <li>- Provide a copy of the annotation guidelines provided to any annotators along with any examples.</li> <li>- Provide any other training or reference material provided to the annotators.</li> </ul>	All	All
	8b	<p><b>If annotation was done, report how many annotators labeled the dataset(s), including the proportion of data in each dataset that were annotated by more than 1 annotator, and the inter-annotator agreement.</b></p> <ul style="list-style-type: none"> <li>- State the number of annotators that were used in total, and what proportion was annotated by multiple individuals</li> <li>- Report the level of agreement among annotators when multiple were used using appropriate statistical measures e.g. Cohen’s kappa.</li> </ul>	All	All
	8c	<p><b>If annotation was done, provide information on the background and experience of the annotators.</b></p> <ul style="list-style-type: none"> <li>- Provide details on the professional background, qualifications, and experience of the annotators.</li> </ul>	All	All

Prompting	9a	<p><b>If research involved prompting LLMs, provide details on the processes used during prompt design, curation, and selection.</b></p> <ul style="list-style-type: none"> <li>- Describe the methodology used to create the initial set of prompts.</li> <li>- Explain the criteria and process used to refine and curate the prompts.</li> <li>- Detail the process used to select the final set of prompts from the curated list.</li> <li>- Describe how prompts were tested to ensure they effectively elicited the desired responses from the LLM.</li> </ul>	All	All
	9b	<p><b>If research involved prompting LLMs, report what data were used to develop the prompts.</b></p> <ul style="list-style-type: none"> <li>- Describe the datasets or sources of information used to create the prompts.</li> <li>- Provide details on the datasets used to evaluate the performance of the prompts.</li> <li>- Report if there was any overlap between the datasets used to develop the prompts and to evaluate the methods.</li> </ul>	All	All
Summarization	10	<p><b>Describe any preprocessing of the data before summarization.</b></p> <ul style="list-style-type: none"> <li>- Outline any preprocessing steps applied to the data before summarization e.g. de-identification.</li> <li>- State if any reformatting or additional processing was performed specifically for summarization e.g. removal of specific sections.</li> </ul>	All	SS
Instruction tuning/Alignment	11	<p><b>If instruction tuning/alignment strategies were used, what were the instructions and interface used for evaluation, and what were the characteristics of the populations doing evaluation?</b></p> <ul style="list-style-type: none"> <li>- Describe the specific instruction/preference datasets provided to the LLM during the tuning or alignment process.</li> <li>- Describe the interface or tools through which evaluators evaluate and provide feedback on the LLM's performance during alignment.</li> <li>- Provide information on the demographics and expertise of the evaluators.</li> </ul>	M D	All
Compute	12	<p><b>Report compute, or proxies thereof (e.g., time on what and how many machines, cost on what and how many machines, inference time, floating-point operations per second (FLOPs)), required to carry out methods.</b></p> <ul style="list-style-type: none"> <li>- Specify the computational resources used, for example machines, time, and cost.</li> <li>- Report the inference time and any metrics related to computational efficiency, as available.</li> <li>- If possible, provide additional metrics such as FLOPs to quantify the computational requirements.</li> </ul>	M D E	All
Ethical Approval	13	<p><b>Name the institutional research board or ethics committee that approved the study and describe the participant-informed consent or the ethics committee waiver of informed consent.</b></p> <ul style="list-style-type: none"> <li>- Name the institutional review board or ethics committee that provided approval.</li> <li>- Describe the informed consent process or the waiver granted by the ethics committee.</li> </ul>	All	All

Open Science	14a	<p><b>Give the source of funding and the role of the funders for the present study.</b></p> <ul style="list-style-type: none"> <li>- Identify the funding sources supporting the research.</li> <li>- Describe any role the funders had in the study design, data collection, analysis, or publication.</li> </ul>	All	All
	14b	<p><b>Declare any conflicts of interest and financial disclosures for all authors.</b></p> <ul style="list-style-type: none"> <li>- Disclose any relationships or activities that could be perceived as influencing the research.</li> <li>- Provide information on any financial interests or affiliations of the authors.</li> </ul>	All	All
	14c	<p><b>Indicate where the study protocol can be accessed or state that a protocol was not prepared.</b></p> <ul style="list-style-type: none"> <li>- Provide details on where and how the clinical study protocol can be accessed by others.</li> </ul>	H	All
	14d	<p><b>Provide registration information for the study, including register name and registration number, or state that the study was not registered.</b></p> <ul style="list-style-type: none"> <li>- If a clinical trial component is undertaken, state the name of the registry and the registration number for the study.</li> <li>- Clearly state if the study was not registered and provide reasons if applicable.</li> </ul>	H	All
	14e	<p><b>Provide details of the availability of the study data.</b></p> <ul style="list-style-type: none"> <li>- Explain where and how the study data can be accessed, including any conditions or restrictions.</li> </ul>	All	All
	14f	<p><b>Provide details of the availability of the code to reproduce the study results.</b></p> <ul style="list-style-type: none"> <li>- Describe how and where the code used in the study can be accessed by others.</li> </ul>	All	All
Public Involvement	15	<p><b>Provide details of any patient and public involvement during the design, conduct, reporting, interpretation, or dissemination of the study or state no involvement.</b></p> <ul style="list-style-type: none"> <li>- Describe how patients or the public were involved in various stages of the research.</li> <li>- Explain if and how the findings were shared with patients or the public.</li> </ul>	H	All
<b>Results</b>				
Participants	16a	<p><b>When using patient/EHR data, describe the flow of text/EHR/patient data through the study, including the number of documents/questions/participants with and without the outcome/label and follow-up time as applicable.</b></p> <ul style="list-style-type: none"> <li>- If EHR data is used, describe the process of how patient/EHR data were selected, filtered, and included in the study.</li> <li>- Specify the number of documents, questions, or participants included and excluded at each stage.</li> <li>- Indicate the number of participants with and without the specific outcome/label and the duration of follow-up.</li> </ul>	E H	All

	16b	<p><b>When using patient/EHR data, report the characteristics overall and, for each data source or setting, and for development/evaluation splits, including the key dates, key characteristics, and sample size.</b></p> <ul style="list-style-type: none"> <li>- Provide a summary of the overall demographic and clinical characteristics of the dataset.</li> <li>- Detail the characteristics for each specific data source or setting.</li> <li>- Describe the sample size and key characteristics for the development and evaluation datasets.</li> </ul>	E H	All
	16c	<p><b>For LLM evaluation that include clinical outcomes, show a comparison of the distribution of important clinical variables that may be associated with the outcome between development and evaluation data, if available.</b></p> <ul style="list-style-type: none"> <li>- Provide a comparison of key predictors, demographics, and clinical characteristics between the development and evaluation datasets.</li> <li>- Report whether the distribution of predictors, demographics, and clinical is comparable between datasets.</li> <li>- These characteristics will depend on the specific context of use and task for each study, as established by literature review and/or domain expert input.</li> </ul>	E H	All
	16d	<p><b>When using patient/EHR data, specify the number of participants and outcome events in each analysis (e.g., for LLM development, hyperparameter tuning, LLM evaluation).</b></p> <ul style="list-style-type: none"> <li>- Report the number of participants and outcome events for each specific analysis.</li> <li>- Describe the stages of analysis and the corresponding data used.</li> </ul>	E H	All
Performance	17	<p><b>Report LLM performance according to pre-specified metrics (see item 7a) and/or human evaluation (see item 7d).</b></p> <ul style="list-style-type: none"> <li>- Report performance overall and for any key subgroups (e.g., sociodemographic, diagnosis, data source).</li> <li>- Consider plots to aid presentation.</li> <li>- Consider reporting confidence intervals overall and for any key subgroups.</li> <li>- Consider reporting uncertainty estimation of the generated output (e.g., LLM-verbalized estimates, logit-based estimates) overall and for any key subgroups.</li> </ul>	All	All
LLM Updating	18	<p><b>If applicable, report the results from any LLM updating, including the updated LLM and subsequent performance.</b></p> <ul style="list-style-type: none"> <li>- Explain any modifications or updates made to the LLM and the reasons behind them.</li> <li>- Report the performance metrics of the updated LLM.</li> </ul>	All	All
<b>Discussion</b>				
Interpretation	19a	<b>Give an overall interpretation of the main results, including issues of fairness in the context of</b>	All	All



		<p><b>the objectives and previous studies.</b></p> <ul style="list-style-type: none"> <li>- Summarize the main findings and their implications overall and for the specified or anticipated healthcare contexts of use.</li> <li>- Discuss any fairness or robustness issues observed, such as biases in predictions.</li> </ul>		
Limitations	19b	<p><b>Discuss any limitations of the study and their effects on any biases, statistical uncertainty, and generalizability.</b></p> <ul style="list-style-type: none"> <li>- Identify and explain the limitations of the study design, robustness of results, and implications for generalisability of findings.</li> </ul>	All	All
Usability of the LLM in context	19c	<p><b>Describe any known challenges in using data for the specified task and domain context with reference to representation, missingness, harmonization, and bias.</b></p> <ul style="list-style-type: none"> <li>- Explain the difficulties encountered in using the data for the specified task e.g. formatting inconsistencies, missingness, class imbalance, or harmonization challenges.</li> <li>- Discuss issues related to data representation and potential biases that may impact findings generalizability or robustness.</li> </ul>	E H	All
	19d	<p><b>Define the intended use for the implementation under evaluation, including the intended input, end-user, level of autonomy/human oversight.</b></p> <ul style="list-style-type: none"> <li>- Specify the purpose of the LLM and the type of input it requires.</li> <li>- Describe the end-users and the level of autonomy or human oversight required.</li> <li>- Discuss barriers to access by the intended end-user, e.g. lack of access to hospital systems with EHRs, wifi, technical support.</li> </ul>	E H	All
	19e	<p><b>If applicable, describe how poor quality or unavailable input data should be assessed and handled when implementing the LLM, i.e., what is the usability of the LLM in the context of current clinical care.</b></p> <ul style="list-style-type: none"> <li>- Explain strategies for managing poor quality or missing input data.</li> <li>- Describe the LLM's usability in real-world clinical settings.</li> </ul>	E H	All
	19f	<p><b>If applicable, specify whether users will be required to interact in the handling of the input data or use of the LLM, and what level of expertise is required of users.</b></p> <ul style="list-style-type: none"> <li>- Describe the extent of user interaction needed for handling input data or operating the LLM.</li> <li>- Specify the level of expertise needed to use the LLM effectively.</li> </ul>	E H	All
	19g	<p><b>Discuss any next steps for future research, with a specific view to applicability and generalizability of the LLM.</b></p> <ul style="list-style-type: none"> <li>- Outline potential areas for further investigation to improve the LLM.</li> <li>- Discuss how the findings can be applied to other contexts or populations</li> </ul>	All	All

LLM = large language model; M = LLM methods; D = *de novo* LLM development; E = LLM evaluation; H = LLM evaluation in healthcare

settings; C = classification; OF = outcome forecasting; QA = long-form question-answering; IR = information retrieval; DG = document generation; SS = summarization and simplification; MT = machine translation, EHR = electronic health record.