npj Digital Medicine 8, 274 (2025)

A framework to assess clinical safety and hallucination rates of LLMs for medical text summarisation

Elham Asgari, Nina Montaña-Brown, Magda Dubois, Saleh Khalil, Jasmine Balloch, Joshua Au Yeung, Dominic Pimenta











Why was this study needed?

As ambient AI scribe tools began rolling out across NHS settings, a major gap remained: how to monitor their safety once in use.

Large language models (LLMs) can generate accurate clinical documentation, but small errors—such as missing negations, misattributions, or slight phrasing shifts—can alter clinical meaning.

No structured framework existed to quantify or mitigate these risks in real-world workflows. This study introduced CREOLA: the first published methodology for proactive, continuous safety monitoring of AI-generated notes in live clinical environments.

What we did

We identified the four most important documentation risks posed by large language models (LLMs): hallucinations, omissions, duplications, and misattributions. Using 450 simulated NHS consultations, we tested how often these errors appeared in AI-generated notes. Each note was reviewed using the CREOLA framework—a tool we designed to flag and classify safety risks. We then tested how changes to system design and prompt structure affected error rates.

The result: a practical method for detecting and reducing risks in real time. CREOLA is now fully embedded in TORTUS and runs automatically on every note.

What we found

Even well-performing LLMs introduce some safety risk:

- On average, 1.5% of note content was hallucinated, and 3.5% of transcript content was omitted
- Of these, 44% of hallucinations and 17% of omissions were rated clinically significant—serious enough to affect care decisions if left uncorrected. But crucially, these risks are fixable. When we refined prompts and outputs, we saw:

- 75% reduction in major hallucinations.
- 58% reduction in minor omission.
- 0 major omissions in the best-performing system configuration.

The refined system performance was equal to or better than known human error rates, demonstrating that with the right safeguards, LLM-generated documentation can meet or exceed accepted standards of clinical accuracy.

What this means for you

CREOLA addresses a critical gap in digital clinical governance: how to oversee LLM safety post-deployment, at scale.

It offers a structured, auditable mechanism to:

- Continuously monitor clinical documentation safety
- Detect emerging risks across specialties
- Support compliance with NHS safety and governance frameworks

TORTUS is the first AVT provider to publish and implement such a system, now integrated into every deployment. For CCIOs, this means real-time safety assurance, traceable audit trails, and alignment with NHS clinical risk standards—marking a new level of confidence in AI-enabled documentation.

Published in partnership with Seoul National University Bundang Hospital



https://doi.org/10.1038/s41746-025-01670-7

A framework to assess clinical safety and hallucination rates of LLMs for medical text summarisation

Check for updates

Elham Asgari^{1,2}⊠, Nina Montaña-Brown¹, Magda Dubois¹, Saleh Khalil¹, Jasmine Balloch¹, Joshua Au Yeung¹ & Dominic Pimenta¹

Integrating large language models (LLMs) into healthcare can enhance workflow efficiency and patient care by automating tasks such as summarising consultations. However, the fidelity between LLM outputs and ground truth information is vital to prevent miscommunication that could lead to compromise in patient safety. We propose a framework comprising (1) an error taxonomy for classifying LLM outputs, (2) an experimental structure for iterative comparisons in our LLM document generation pipeline, (3) a clinical safety framework to evaluate the harms of errors, and (4) a graphical user interface, CREOLA, to facilitate these processes. Our clinical error metrics were derived from 18 experimental configurations involving LLMs for clinical note generation, consisting of 12,999 clinician-annotated sentences. We observed a 1.47% hallucination rate and a 3.45% omission rate. By refining prompts and workflows, we successfully reduced major errors below previously reported human note-taking rates, highlighting the framework's potential for safer clinical documentation.

One of the most appealing applications of LLMs in healthcare is for administrative tasks¹. Clinicians devote a substantial amount of time to documentation², and prolonged interaction with electronic health records, where clinical documentation is logged, has been demonstrated to raise cognitive load and lead to burnout³. In fact, the use of LLMs for clinical documentation, especially clinical note generation⁴ or consultation summarisation^{5,6}, is an active area of research.

However, LLMs are known to produce errors in many settings, from document summarisation⁷, to general reasoning tasks as well as more clinically relevant tasks⁸. These errors can be categorised as "hallucinations": known as an event where LLMs generate information that is not present in the input data, or omissions: the event where LLMs miss relevant information from the original document. Errors in clinical documentation generation can lead to inaccurate recording and communication of facts^{10,11}. Inaccuracies in the document summarisation task can introduce misleading details⁸ into transcribed conversations or summaries, potentially delaying diagnoses¹² and causing unnecessary patient anxiety.

The problem of hallucinations poses a significant challenge to date^{1,13}. The occurrence of hallucinations has previously been attributed to the data quality during model training^{14,15}, the type of model training methodology¹⁶ and prompting strategies¹⁷.

Recent work has established that hallucination may be an intrinsic, theoretical property of all LLMs⁹. Consequently, there is a growing body of work focused on the technical evaluation of LLM accuracy and the detection

and mitigation of hallucinations in LLMs¹⁸. However, the prevalence, causation, and evaluation of hallucinations in a clinical context, as well as their subsequent impact on clinical safety, remains an open question.

Clinical documentation can be variable in quality^{19,20}, and studies estimate that human-generated clinical notes have, on average, at least 1 error and 4 omissions²¹. Given the increased usage of LLMs for clinical documentation^{22,23}, several methods have been proposed for evaluating clinical documentation generated using LLMs.

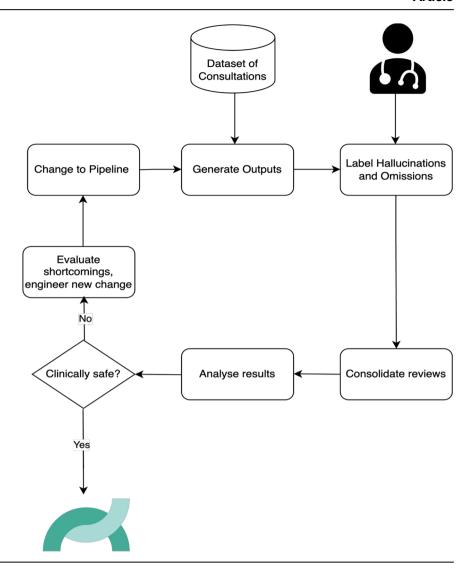
Relevant clinical evaluation frameworks typically include categorising clinical errors for downstream analysis. Typically, these differ from traditional natural language processing (NLP) taxonomies¹⁶, which have separated hallucination types into distinct categories, for example, into "intrinsic" and "extrinsic"²⁴, "factuality" and "faithfulness"¹⁶, "factual mirage" and "silver lining"²⁵ errors. The differences between general and clinical taxonomies arise from the necessity of increased granularity of clinical error types, which are not captured by the broader, general methods.

For example, Tierney et al.²⁶ propose using a modified version of the Physician Documentation Quality Instrument-9, accounting for hallucinations and bias, while Abacha et al.²³ propose evaluating clinical note quality using automated metrics. However, these relevant clinical categorisations have not assessed the implications of the mistakes for downstream harm.

Automated metrics, such as Recall-oriented Understudy for Gisting Evaluation (ROUGE)²⁷, Bilingual Evaluation Understudy (BLEU)²⁸ and

¹Tortus AI, London, UK. ²Guy's and St Thomas NHS Trust, London, UK. 🖂 e-mail: asgelham@gmail.com

Fig. 1 | Our workflow for the assessment of LLM output using CREOLA platform. This diagram illustrates the process we followed using the available dataset for various experiments, including the input from clinicians for labelling and consolidating reviews, followed by a safety analysis.



Bidirectional Encoder Representations from Transformers (BERT)-score²⁹, while useful for comparing model-generated text with expert-written examples, exhibit significant limitations when applied to the evaluation of healthcare-related content. These metrics, primarily focused on surface-level textual similarity, fail to capture the semantic nuances, contextual dependencies, and domain-specific knowledge crucial for accurate medical discourse³⁰. This deficiency is particularly problematic in healthcare settings, where understanding complex medical concepts (e.g., symptoms, diagnoses, treatments) and their interrelationships is paramount for patient well-being and effective decision-making.

Despite the exponential growth in benchmarks for model reasoning abilities³¹, the evaluation of LLMs on clinical tasks has typically been carried out via "question-answering" (QA) benchmarks^{5,8,32}. These tasks assess models' accuracy over various clinical questions, typically derived from licensing exams. While these methods offer insights into the factual knowledge and reasoning abilities of LLMs, they do not assess clinical or medical capabilities such as medical text summarisation.

Singhal et al.³³ have outlined the challenges of evaluating LLMs in various medical contexts, including medical exams, research and consumer queries. They have proposed a human evaluation model for the answers provided by different LLMs that checks on factuality, precision, possible harm and bias. Other evaluation factors such as fairness, transparency, trustworthiness and accountability have been suggested in using LLMs in healthcare³⁴. In a more recent study, Tang et al. assessed human evaluation based on metrics such as coherence, factual consistency, comprehensiveness and potential harm.

Interestingly, they assessed the clinician's preference for different outputs³⁵. More recently, Tam et al.³⁶ have introduced QUEST as a framework for human evaluation of LLMs in healthcare following a comprehensive literature review on the topic. QUEST includes five principles for human evaluation of LLMs, including Quality of information, Understanding and reasoning, Expression style and persona, Safety and harm, and Trust and confidence.

Multiple benchmarks have been proposed to evaluate model summarisation capabilities in the biomedical domain, including over biomedical literature^{37–40}, medical forum conversations⁴¹, and radiology reports^{22,42}. However, these benchmarks do not capture the nuances of patient-facing clinical interactions, where LLM-documentation holds most promise.

Recently, Umapathi et al.⁴³ have assessed models' tendency towards hallucination. They reported that LLMs were significantly variable in their accuracy depending on the prompts used. However, the MedHALT benchmark is limited to assessing LLM's reasoning capabilities over the medical domain in a QA format. Most relevantly, Moramarco et al.²¹ benchmark BART models on the PriMock dataset and find that they produce 3.9 errors and 6.6 omissions on average per note. However, they did not assess the model's impact or human errors on patient safety as part of their study.

This study aims to contribute to the ongoing effort to ensure clinical safety in using LLMs for note generation by introducing a framework which has four components: (1) a clinically and technically-informed error taxonomy to classify LLM outputs, (2) an experiment structure to comprehensively and iteratively compare outputs within our LLM document generation pipeline, (3) a clinical safety framework to assess potential harms

of errors in LLM outputs, and (4) an encompassing graphical user interface (GUI), CREOLA, to perform and assess all previous steps. Figure 1 shows our workflow based on the framework. We present our findings and insights from applying our framework, which, to our knowledge, represents the largest manual evaluation of LLM clinical note generation to date.

Our objective is to promote the efficient, reliable, and confident use of LLMs for clinical documentation, thus supporting healthcare providers in

delivering high-quality care and overall reducing the administrative workload for clinicians.

Results

Dataset

We conducted a series of 18 experiments, each consisting of 25 primary care consultation transcripts from the PriMock dataset⁴⁴ For each transcript, we

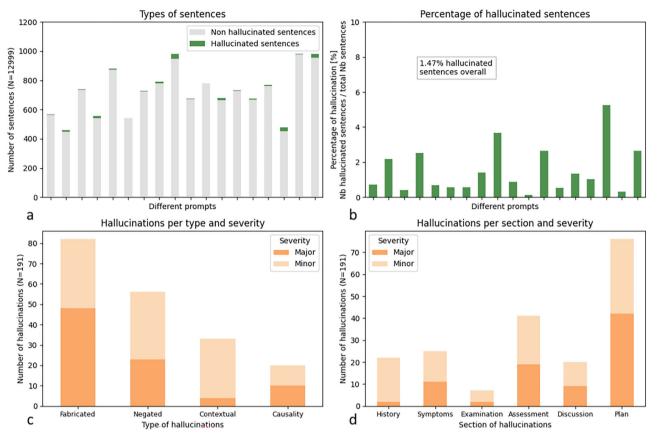
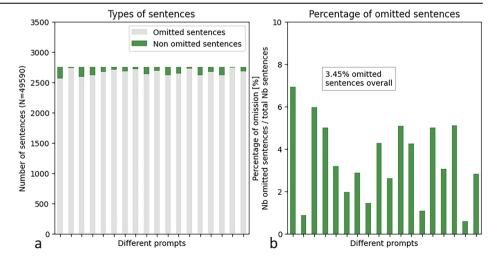


Fig. 2 | Incidence of hallucinations, their types, the section of the note they appear in, and clinical risk. The figure illustrates the occurrence of hallucinations in the generated sentences based on different prompts (a) and their corresponding

percentages (b). It also highlights the type and clinical severity of hallucinations (c) and the specific sections of the note where they appeared (\mathbf{d}).

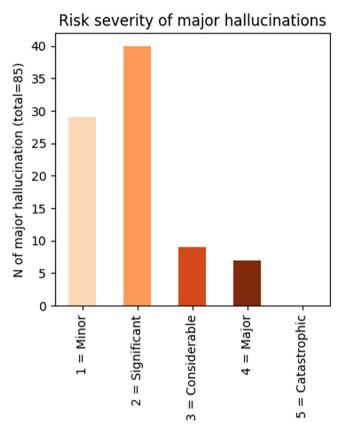
Fig. 3 | Number of omitted sentences based on different prompts. The figure illustrates the number of omitted sentences using different prompts (a) and their respective percentage (b).



generated paired clinical documentation using an LLM, resulting in 450 consultation transcript-note pairs. This was a total of 49,590 transcript sentences and 12,999 clinical note sentences that were manually evaluated and labelled for any hallucinations or omissions.

Experiments

Our experiments were guided by our framework, which provides a systematic approach to evaluate LLM outputs quantitatively. Using a baseline LLM prompt and workflow, we generated 25 transcript-note pairs. We recruited 50 medical doctors for manual evaluation. For each paired transcript-note, we had two clinician reviewers evaluate each sentence in the



 $\label{lem:Fig. 4 | Severity of risk in major hallucinations. We assessed the clinical risk resulting from major hallucinations based on our suggested framework.}$

clinical note to ensure it was evidenced in the transcript; sentences that were not evidenced were labelled as hallucinations. We also highlighted each sentence in the transcript for review, checking if it was present in the output note if it was clinically relevant; and if not, they were labelled as an omission. If the hallucination or omission could change the diagnosis or management of the patient (if left uncorrected), it was marked as 'Major', otherwise, they would be labelled as 'minor'. In cases of discrepancy between the two reviewers, consolidation was performed by a senior clinician with over 20 years of clinical experience.

We additionally identified the specific sections of the notes where the hallucinations occur (main history, examination, discussion, symptoms assessment, and plan). The result of each experiment informed our subsequent experiment approach to analyse how prompt/ workflow/ engineering changes affect the hallucinatory potential for clinical note generation.

All experiments were conducted on CREOLA, our in-house platform, designed to allow clinicians to identify and label relevant hallucinations and omissions in clinical text. Using this platform, we were able to implement our framework to quantify and track changes in our prompts and model configurations and iteratively modify our approach to ensure the safe integration of LLM-generated summaries into clinical practice.

Hallucinations

Of 12,999 sentences in 450 clinical notes, 191 sentences had hallucinations (1.47%), of which 84 sentences (44%) were major (could impact patient diagnosis and management if left uncorrected). Of the hallucination types, 82 (43%) were fabricated, 56 (30%) were negations, 33 (17%) were contextual, and 20 (10%) were related to causality.

Major hallucinations occurred in all sections, but most commonly in Plan (21%), Assessment (10.5%), and Symptoms (5.2%) sections. The most common hallucination type were fabrications and primarily appeared in the planning section of the clinic note (Fig. 2). Examples of the various hallucination types are available in supplementary materials.

Omissions

Of the 49,590 sentences from our consultation transcripts, 1712 sentences were omitted (3.45%), of which 286 (16.7%) of which were classified as major and 1426 (83.3%) as minor. Figure 3 shows the number of omissions and the percentage of omitted sentences based on different prompts.

Grading hallucinations and omissions by clinical safety impact

Inspired by protocols in medical device certifications, we applied the clinical safety assessment of our framework described in the methods section. We classified the clinical risk (Major or Minor) and evaluated the risk severity of all identified major hallucinations as depicted in Fig. 4. We also determined

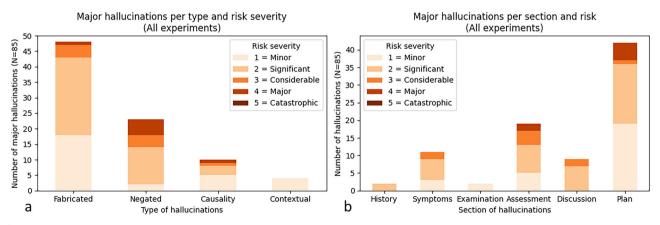
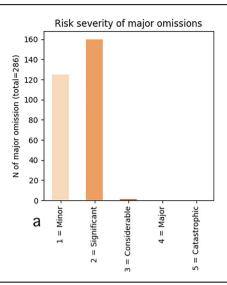


Fig. 5 | Hallucination risk assessment. We have categorised the clinician severity of hallucination risks based on their type (a) and the section of the clinical notes where they occurred (b).

Fig. 6 | Severity of clinical risk for major omissions and the section of the clinical note where they most occurred. This figure illustrates the clinical risk severity due to major omissions (a) and indicates the sections of the notes where they are most likely to occur (b).



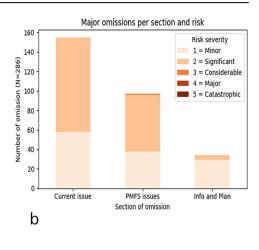
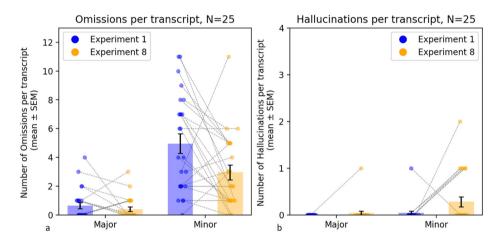


Fig. 7 | A comparison of omissions and hallucinations between Experiment 1 and Experiment 8. Comparing omissions (a) and hallucinations (b) using the base prompt from Experiment 1 and the updated style prompt from Experiment 8. Each connected dot represents the change in hallucinations or omissions for a given document.



the risk of hallucinations based on their type and where in the sentence they occurred (Fig. 5).

We assessed the clinical risk resulting from major hallucinations based on our suggested framework inspired by protocols in medical device certifications.

We conducted the risk assessment on the omissions (Fig. 6a) and detailed the note section that would likely have been affected (Fig. 6b). Major omissions were most common in current issues, followed by PMFS, and Info and Plan sections (55%, 35%, and 10%, respectively).

Examples of hallucinations and omissions are available in the supplementary materials.

Iterative experiments can significantly reduce hallucination and omission rates in LLM-generated clinical notes

Through a series of 18 iterative experiments, we tested a combination of prompting and workflow strategies, including structured prompting, atomisation, function calls and JSON-based outputs, an additional LLM revision step, and templating (SOAP -Subjective, Objective, Assessment, plannote), which are explained in more detail in the methods section.

Modifying the prompt from the baseline used in Experiment 1 to include a style update used in Experiment 8, resulted in a reduction of both major and minor omissions. Although there was a slight increase in hallucinations in Experiment 8, these were mostly minor. Figure 7 illustrates the number of hallucinations and omissions recorded in Experiments 1 and 8.

We then compared the outputs using various structured prompts in Experiments 3 and 8 illustrated in Table 1.

The change in the prompt from Experiments 3 to 8 reduced the incidence of major hallucinations by 75% (from 4 to 1), major omissions by 58% (24 to 10), and minor omissions by 35% (114 to 74) (Fig. 8). By following structured prompting, including a style update and instructing the model to output the status "unknown" for instances where information was missing from the transcript, we significantly improved performance.

In a subsequent experiment (Experiment 5), we found that incorporating a chain-of-thought prompt (Table 1, supplementary material), to extract facts from the transcript—a process referred to as atomisation—before generating the clinical note, led to an increase in major hallucinations and omissions. Figure 9 shows the comparison between omissions and hallucinations between this experiment and our base Experiment (Experiment 1).

We also compared the output of Experiment 5 to Experiment 3, which used structured prompts. We found that major hallucinations increased from 4 to 25, minor hallucinations from 5 to 29, major omissions from 24 to 47, and minor omissions from 114 to 188 (Fig. 10). This result precluded the new change from being evaluated for clinical safety, as the increase in hallucinations and omissions was considered too large to be considered useful

We found that using prompts with function calling was useful in ensuring that the outputs adhered to a specific structure required for different electronic health records. Utilising our framework, we iteratively

Table 1 | Prompt changes that led to decreased hallucinations and omissions

Experiment 3 Experiment 8 You are a medical office assistant drafting documentation for a physician. DO NOT ADD any content You are a highly accurate medical office assistant drafting that isn't specifically mentioned IN THE TRANSCRIPT. From the attached transcript generate a SOAP documentation for a physician. Every decision you take is life or death and must be 100% note based on the below template format for the physician to review, include all the relevant information accurate. DO NOT ADD any and do not include any information that isn't explicitly mentioned in the transcript. If nothing is content that isn't specifically mentioned IN THE TRANSCRIPT. mentioned just return [NOT MENTIONED]. From the attached transcript It is VITAL that all the information in the note is as accurate as possible. Avoid repeating the same generate a clinical note based on the below template format for information in different sections where possible. Write the note from the perspective of the physician. the physician to review, Only include any section of the template if there is information from the transcript, otherwise omit it. include all the relevant information and do not include any information that isn't explicitly mentioned in the transcript. If nothing is mentioned just return INOT MENTIONEDI. It is vital that all the information in the note is as accurate as possible. Avoid repeating the same information in different sections where possible. Write the note from the perspective of the physician. DO NOT add associate or relate causes for medical conditions unless explicitly specified by the Physician. See below for a template to outline the structure of the output and style preferences to follow. Template: Template for Clinical SOAP Note Format: Referral Reason/reason for appointment: - HPI: [include here any mentioned symptoms, chronological narrative of patients complaints, History information obtained from other sources(always identify source if not the patient).] - Allergies - Medications - Past medical history. [include here all of the patients past conditions, treatments and encounters, also - History of presenting complaint include relevant social history here including smoking, alcohol, drug use and occupation/travel - Past Medical History history] - Family/Social History - Review of systems [include here any additioinal symptoms in other organs that is relevant to the initial presentation] Sensitive information Observations: - Current medications [list medicines out each on a seperate line, in a standard format where the - Examination findings information is mentioned: [DRUG NAME][DRUG DOSE][DRUG FREQUENCY][INDICATION] - Investigation results Objective: - Impression or clinical assessment - Vital signs [including any mentioned blood pressure, pulse rate, oxygen saturation, temperature] Plan. - Physical exam [the examination findings from the physical exam, if mentioned] - Planned investigations - Test Results [include in this section any lab test results or imaging reports] Assessment / Problem List: - Assessment: [A one sentence description of the patient and major problem as described >by the - New prescribed medication or therapies physician, including the diagnosis the physician has identified] - Problem list: [A numerical list of clinical problems arising from this encounter and active ongoing - Communication, reassurance & patient understanding of care - Actions for referrer/GP medical problems the patient has. Present each problem as [Condition][Status:active/suspected/ confirmed/past/unknown], list each problem on a separate line, leave status as unknown if not - Explained medical terms mentioned in the transcript] Plan: [include here any management plan mentioned in the transcript, including patient education, prescriptions, tests, referrals or other plans.] Follow-up: [include here any plan mentioned to see the patient again, or to be discharged.] Style preferences: Please adhere to the following style guidelines: - Write from the perspective of the physician (first person) - Write from the perspective of the physician (first person) - Be ultra concise - Be ultra-concise - Be ultra-precise, do not use generalising terms - Use bullet points and broken sentences - Be highly detailed - Include ALL important negations in the relevant sections (e.g. the patient has no fever) the clinician has

- elicited as well as all positive findings.
- Use bullet points and single words, not sentences.
- Always list medications in a list in the following format for each one: medicine, dose, frequency, indication
- Always document if drug allergies are present or not
- Examination findings always refer to a physical exam, only include signs here, not symptoms
- Preserve quantities if mentioned in the text

improved the performance of the structured notes across several experiments (6, 9, 10, and 11). From the first to the last iteration (Experiments 6 to 11), we made meaningful improvements to the prompts, including instructions on adherence to subheadings and the addition of a writing style guidance (e.g., a list of writing rules to follow). Table 2 shows the prompts used in the four function call experiments.

As a result of these changes, we eliminated major omissions completely, decreasing them from 61 to 0, and reduced minor omissions by 58%, from 130 to 54. Additionally, we lowered the total number of hallucinations by 25%, reducing them from 4 to 3 (Fig. 11).

We then examined the two best-performing experiments with the fewest hallucinations and omissions (Experiments 8 and 11) for the type of hallucinations they produced and where they were more likely to appear in the sentence (Fig. 12).

Experiment 11 did not have any major omissions, and the risk assessment of major omissions for Experiment 8 is shown in Fig. 13.

Discussion

Our study supported that hallucinations and omissions may be intrinsic theoretical properties of current LLMs⁹. LLMs can output unfactual or unfaithful text with high degrees of confidence⁴⁵ which can be particularly dangerous in a high-stakes environment such as healthcare. Our framework quantifies the clinical impact and implications that LLM omissions and hallucinations may lead to if unchecked or uncorrected; only then can

clinical safety be meaningfully addressed. Once LLM errors are identified and quantified, we can make iterations on LLM prompts, workflows or engineering design to reduce or eliminate these errors. By concentrating on minimising errors that could significantly impact patient care, we can align LLMs with clinical safety standards and regulations.

To our knowledge, we have conducted the largest manual evaluation on the task of LLM clinical note generation to date. To implement our framework at scale we built CREOLA, an in-house platform, designed to enable clinicians to identify and label relevant hallucinations and omissions in clinical text and inform future experiments. However, similar publicly available platforms (e.g. labelbox, explosion AI, autoblocks etc.) may be used for this task. Experiments within CREOLA can validate or discredit architectures and prompt approaches in a safe sandbox environment before clinical deployment.

Our experiment results show that omissions are more common than hallucinations (3.45% to 1.47%, respectively). However, hallucinations were much more likely to be classified as a "major" error compared to omissions (44% to 16.7%, respectively). This means that hallucinations are more likely to lead to downstream harm and impact clinical care if left uncorrected. Hallucinations occurred in all sections, but most commonly (20%) in the Plan section of the clinical note. This is an important finding as this section often contains direct instructions or actions to colleagues or patients that can impact clinical safety. The most concerning hallucinations were the negation type (30% of total hallucinations). These mostly appeared in the planning section and contradicted what was said during the consultation. Negation hallucinations can lead to significant confusion and harm, and without the full context of the consultation, readers may struggle to discern which negation is true and which is false.

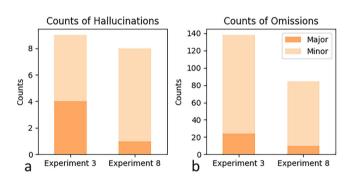


Fig. 8 | Comparison of hallucination and omission counts between two experiments, assessing differences in prompt engineering effect on the quality of outputs. The figure shows the changes in the counts of hallucinations (a) and omissions (b) between Experiments 3 and 8 which have resulted from the changes in the prompts.

By designing prompts that addressed specific aspects of the notes (base, template and style), we were able to focus our iterations to achieve the best results. However, it is important to note that summarisation tasks require the ablation of certain data from the original text to make it a concise, relevant, and useful summary artefact. Optimal omission rates depend significantly on the context, the quality of input and the reviewer receiving output. Our framework allows us to focus on the clinical impact of omissions and hallucinations quantitatively. Overall, our hallucination rates are similar to those reported in the literature for generalist tasks⁴⁶. For the clinical summarisation task, Experiment 8 achieved 1 major hallucination and 10 major omissions, whilst experiment 11 achieved 2 major hallucinations and 0 major omissions over 25 notes. These results are highly encouraging, as our iterative experiment process has resulted in fewer errors per note than those reported in the clinical literature. Moramarco et al.²¹ reported 3.9 errors and 6.6 omissions per note as produced by a BART model and 1 error and 4 omissions per human-written note. This improvement is likely due to the large parameter sizes of modern large language models in combination with our framework. Although this rate is subject to change depending on the text and experiment, our results suggest that we can achieve state-of-the-art, sub-human clinical error rates by carefully engineering and subsequently validating LLMs to produce safe outputs.

Our study is limited in several ways. Firstly, the sample size of medical transcripts used was relatively small; the sample size was chosen to balance the trade-off of annotation volume required for the comparison of different experiments against sample size and number of experiments performed. Additionally, we only evaluated one LLM (GPT-4), selected due to its

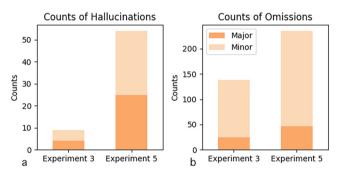


Fig. 10 | Comparison of hallucination and omission counts between two experiments, assessing the efficacy of a data-extraction intermediate step versus a normal note-generation step. The figure compares the number of hallucinations (a) and omissions (b) between Experiment 3, where the standard note generation step was used, and experiment 5, where data extraction was performed before note generation.

Fig. 9 | A comparison of omissions and hallucinations between Experiment 1 and Experiment 5. Comparing omissions (a) and hallucinations (b) in our base experiment (Experiment 1) with Experiment 5 where the transcript was broken down into concise facts before being passed to the LLM for note generation. Each connected dot represents the change in hallucinations or omissions for a given document.

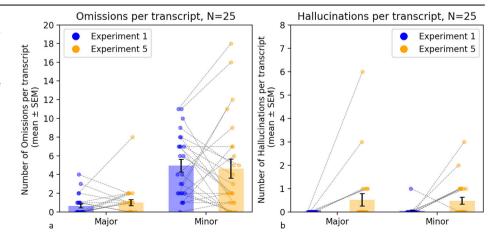


Table 2 | Prompts used in experiments 6, 9, 10 and 11

Eperiment 6

Please use the function below to generate a customised output. provide your output in json format.break

For each parameter value you provide, make sure to include all properties defined in the schema. If a parameter is an array please try to separate ideas into separate items

You are a highly accurate medical officer drafting documentation for a physician. You will receive a transcript of a medical consultation between a patient and a clinician. Your task is to identify the patient's key problems in that encounter, and then extract information provided to you in the format outlined in a JSON schema for each individual problem. A problem is a single discrete issue for the patient, encompassing presenting complaint, associated symptoms, and relevant history e.g. shortness of breath, needs new housing, recent bereavement etc.

It is VITAL that you include all properties mentioned in the schema, if there is a field that is not mentioned in the transcript just write [NOT MENTIONED]. Only include information strictly mentioned in the transcript. Failure to do so may cause harm to the patient. DO NOT duplicate information in more than one problem.

Aim to capture all problems mentioned in the transcript into their own array items.

Experiment 9

Please use the function below to generate a customised output. provide your output in ison format.

For each parameter value you provide, make sure to include all properties defined in the schema. If a parameter is an array please try to separate ideas into separate items

You are a highly accurate medical officer drafting documentation for a physician. You will receive a transcript of a medical consultation between a patient and a clinician. Your task is to identify the patient's key problems in that encounter, and then extract information provided to you in the format outlined in a JSON schema for each individual problem. A problem is a single discrete issue for the patient, encompassing presenting complaint, associated symptoms, and relevant historye.g. shortness of breath, needs new housing, recent bereavement etc.

It is VITAL that you include all properties mentioned in the schema, if there is a field that is not mentioned in the transcript just write [NOT MENTIONED]. Only include information strictly mentioned in the transcript. Failure to do so may cause harm to the patient. DO NOT duplicate information in more than one problem.

For each parameter you provide in the tool call, please adhere to the following style guidelines:

- Write from the perspective of the physician (first person)
- Be ultra-concise
- Be ultra-precise, do not use generalising terms
- Be highly detailed
- Include ALL important negations in the relevant sections (e.g. the patient has no fever) the clinician has elicited as well as all positive findings
- Use bullet points and single words, not sentences
- Always list medications in a list in the following format for each one: medicine, dose, frequency, indication
- Always document if drug allergies are present or not
- Examination findings always refer to a physical exam, only include signs here, not symptoms
- Preserve quantities if mentioned in the text
- Avoid repeating the same information in different sections where possible

Experiment 10

Please use the function below to generate a customised output. provide your output in ison format.

For each parameter value you provide, make sure to include all properties defined in the schema. If a parameter is an array please try to separate ideas into separate items

You are a highly accurate medical officer drafting documentation for a physician. You will receive a transcript of a medical consultation between a patient and a clinician. Your task is to identify the patients key problems in that encounter, and then extract information provided to you in the format outlined in a JSON schema for each individual problem. A problem is a single discrete issue for the patient, encompassing presenting complaint, associated symptoms, and relevant historyeg. shortness of breath, needs new housing, recent bereavement etc.

It is VITAL that you include all properties mentioned in the schema, if there is a field that is not mentioned in the transcript just write [NOT MENTIONED]. Only include information strictly mentioned in the transcript. Failure to do so may cause harm to the patient. DO NOT duplicate information in more than one problem.

For each parameter you provide in the tool call, please adhere to the following style guidelines:

- Write from the perspective of the physician (first person)
- Be ultra-concise
- Be ultra-precise, do not use generalising terms
- Be highly detailed
- Include ALL important negations in the relevant sections (e.g. the patient has no fever) the clinician has elicited as well as all positive findings
- Use bullet points and single words, not sentences
- Always list medications in a list in the following format for each one: medicine, dose, frequency, indication
- Always document if drug allergies are present or not
- Examination findings always refer to a physical exam, only include signs here, not symptoms
- Preserve quantities if mentioned in the text
- Avoid repeating the same information in different sections where possible

Experiment 11

Please use the function below to generate a customised output. provide your output in JSON format

For each parameter value you provide, make sure to include all properties defined in the schema. If a parameter is an array please try to separate ideas into separate items.

You are a highly accurate medical officer drafting documentation for a physician. You will receive a transcript of a medical consultation

between a patient and a clinician. Your task is to identify the patients key clinical problems (the presenting complaint or complaints) in

that encounter, and then extract information provided to you in the format outlined in a JSON schema for each individual problem. A

problem is a single discrete issue for the patient, encompassing presenting complaint, associated symptoms, and relevant history -

e.g. shortness of breath, needs new housing, recent bereavement etc.

It is VITAL that you include all properties mentioned in the schema, if there is a field that is not mentioned in the transcript just write

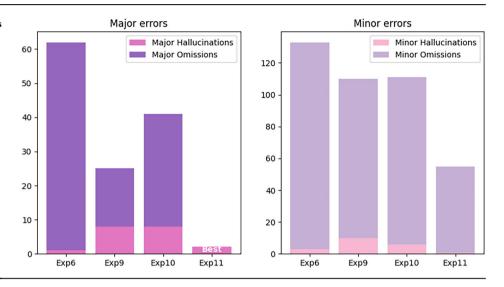
[NOT MENTIONED]. Only include information strictly mentioned in the transcript. Failure to do so may cause harm to the patient. DO

NOT duplicate information in more than one problem.

For each parameter you provide in the tool call, please adhere to the following style guidelines:

- Write from the perspective of the physician (first person)
- Be ultra-concise
- Be ultra-precise, do not use generalising terms
- Be highly detailed
- Include ALL important negations in the relevant sections (e.g. the patient has no fever) the clinician has elicited as well as all positive findings.
- Use bullet points and single words, not sentences
- Always list medications in a list in the following format for each one: medicine, dose, frequency, indication
- Always document if drug allergies are present or not
- Examination findings always refer to a physical exam, only include signs here, not symptoms
- Preserve quantities if mentioned in the text
- Avoid repeating the same information in different sections where possible.
- Ensure the output only has the headings Description of Problem, History,
 Examination and Comments and no other headings

Fig. 11 | Number of hallucinations and omissions in the function call experiments. This figure illustrates how the change in prompts within the function call experiments can change the number of major and minor hallucinations and omissions in the output.



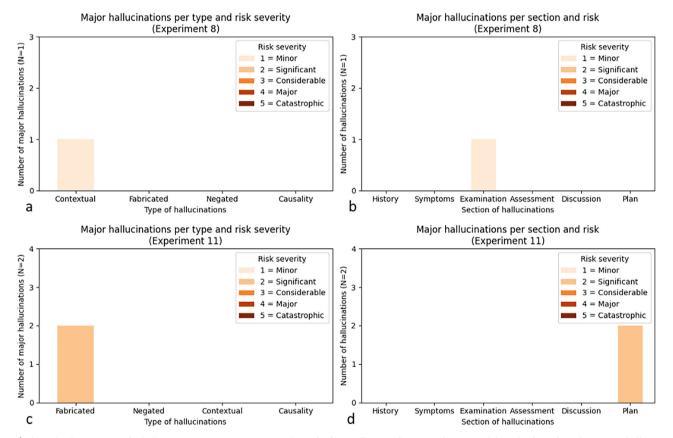


Fig. 12 | Clinical risk assessment for hallucinations in experiments 8 and 11. The figures illustrate the type and severity of clinical risk resulting from major hallucinations in Experiments 8 (a) and 11 (c) and the section of the notes they occur (b and d).

established performance in text summarisation at the time of our experiments. The ongoing development and enhancement of open-source large language models (LLMs) are likely to boost their application in medicine⁴⁷. Recently, Zhang et al. reviewed how fine-tuning open-source LLMs such as PRIMERA, LongT5, and Llama-2 can enhance their ability to summarise medical evidence effectively⁴⁸. Furthermore, our experiments use a direct prompting scheme (Supplementary Data 1). Newer methods such as (but not limited to) Retrieval-Augmented Generation (RAG)⁴⁹, Chain of Thought (CoT)⁵⁰, or the use of knowledge graphs⁵¹ have recently been used to enhance the performance of LLMs. For example, by equipping LLMs with

domain-specific knowledge, RAG enables the models to generate more precise and pertinent results^{52,53}, whilst CoT generally enhances model reasoning abilities. A straightforward extension of this work is using this framework over different experimental configurations, such as using different models or prompting techniques, and comparing the impact on reported performance to clinical safety metrics.

Finally, using human annotators to evaluate large amounts of data is expensive and unsustainable. In the long run, the automated evaluation of model output⁵⁴ is a consequential future direction which will enable the scalable assessment of a larger volume of information, with clinicians

remaining in the loop by "supervising" evaluator models via the inspection of a sub-sample of the outputs. 'LLM-as-a-Judge', which refers to an LLM tasked with evaluating, scoring, or assessing the quality, correctness, or appropriateness of outputs, has recently been described⁵⁵ and its applicability has been discussed⁵⁶. Utilising the capabilities of LLMs for initial screening can significantly reduce the time and resource demands on human evaluators. This can make the CREOLA framework more scalable

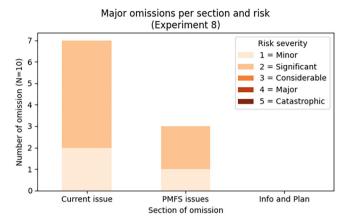


Fig. 13 | Assessment of clinical risk resulting from major omissions. The figure illustrates the clinical risk from major omissions in Experiment 8 and the section of the clinical notes where it has occurred.

Table 3 | The likelihood of a hazard occurring

Likelihood Category	Interpretation
Very high	Certain or almost certain; highly likely to occur
High	Not certain but very possible: reasonably expected to occur in the majority of cases
Medium	Possible
Low	Could occur but in the great majority of occasions will not
Very low	Negligible or nearly negligible possibility of occurring

and efficient over time. This hybrid approach aligns with ongoing advancements in AI and has the potential to maintain rigorous oversight while ensuring scalability.

In this work, we present a framework for the clinical safety assessment of LLMs in clinical documentation scenarios. Using the CREOLA platform, we analyse the impact of prompting techniques on the safety of LLM outputs. Our iterative modification process allows us to reach new low hallucination and omission rates - our best-performing experiments outperform previously reported model and human error rates - facilitating confident deployment of our solutions to end clinical users. Additionally, CREOLA provides a sandbox environment which buffers users and patients from harm in case the iteration leads to higher clinical error rates. The addition of clinical safety assessment to prompt evaluation creates a valuable framework for implementing note summarisation tools in clinical practice. We propose that our suggested framework, which combines the assessment of hallucinations and omissions with an evaluation of their impact on clinical safety, can serve as a governance and clinical safety assessment template for various organisations. This approach aims to empower clinicians to become key stakeholders in the deployment of large LLMs in clinical settings.

Methods

We propose a multi-component framework to evaluate hallucinations and omissions in clinical documentation generated by LLMs. Central to our approach is the concept of "clinician in the loop". Given their expertise, clinicians are uniquely positioned to identify clinical errors made by the models, making their involvement essential. A specialised annotation platform (CREOLA) was developed to facilitate clinician labelling for each experimental dataset.

Experimental design: Our experiments systematically assessed how various prompting techniques and workflow structures influenced the accuracy and reliability of clinical notes derived from primary care consultation transcripts. Typical parameters varied in our experiments included the complexity and specificity of prompts (such as the addition of structured sections, negations, or perspective changes), as well as the number of LLM calls, such as introducing an additional revision step through a secondary LLM call. For consistency and reproducibility, we used OpenAI's GPT-4 (GPT-4-32k-0613), setting the seed to 210, temperature to 0, and a top-p value of 0.95 to accommodate clinical language complexities.

Table 4 | Guidance for assessing the level of harm

Consequence classification	Interpretation	Number of patients affected
Catastrophic	Death	Multiple
	Permanent life-changing incapacity and any condition for which the prognosis is death or permanent life-changing incapacity, severe injury or severe incapacity from which recovery is not expected in the short term	Multiple
Major	Death	Single
	Permanent life-changing incapacity and any condition for which the prognosis is death or permanent life-changing incapacity, severe injury or severe incapacity from which recovery is not expected in the short term	Single
	Severe injury or severe incapacity from which recovery is expected in the short term	Multiple
	Severe psychological trauma	Multiple
Considerable	Severe injury or severe incapacity from which recovery is expected in the short term	Single
	Severe psychological trauma	Single
	Minor injury or injuries from which recovery is not expected in the short term	Multiple
	Significant psychological trauma	Multiple
Significant	Minor injury or injuries from which recovery is not expected in the short term	Single
	Significant psychological trauma	Single
	Minor injury from which recovery is expected in the short term	Multiple
	Minor psychological upset; inconvenience	Multiple
Minor	Minor injury from which recovery is expected in the short term; Minor psychological upset; inconvenience; any negligible severity	Single

Table 5 | Calculating the likelihood of an error occurring in the text output

	Per 25 examples	Possibility
Very High	22.5	90%
High	15	60%
Medium	7	10-60%
Low	2.5	10%
Very Low	0.5	1%

Likelihood	Very high	3	4	4	5	5
	High	2	3	3	4	5
	Medium	2	2	3	3	4
	Low	1	2	2	3	4
	Very low	1	1	2	2	3
		Minor	Significant	Considerable	Major	Catastrophic
		Conse	quence			

Fig. 14 | Risk estimation based on the likelihood and consequence of harm occurrence. This figure illustrates the scoring of clinical risk based on the likelihood of an incident occurring and the severity of harm it may cause.

All prompts used are detailed in Supplementary Materials, Table 4. To achieve a meaningful clinical comparison of efficacy and safety in a data-driven way, our framework relies on the definition of a 'baseline' experiment against which to compare results. The baseline experiment must use the same input data points as the new experiment. To clearly attribute an experiment's results to a specific change, we aim only to alter one parameter from the baseline experiment configuration at a time.

We used different methodologies to assess and improve model output as described below:

Model Improvements: This outlines modifications to individual LLM calls within our workflow while preserving the same overall structure. Common modifications include the prompt, the model used for the call, or the model hyperparameters, such as maximum output tokens (Table 1).

Workflow Improvements: We implemented changes to explore new methods for generating a specific type of output. For instance, in our clinical note generation based on a transcript, we decided to extract a list of facts from the transcript before making a single call to the LLM for the final note (Supplementary Table 1). We then evaluated how this approach affected the frequency of hallucinations and omissions (Figs. 9 and 10). Additionally, we included an extra LLM call in some experiments to improve the quality of the output (Supplementary Fig. 1).

Clinician vs LLM generated notes: Several members of our clinical team were tasked with creating notes based on consultation transcripts. We then utilised the framework to identify any hallucinations and omissions in these notes, which allowed us to compare the clinician-created notes with those generated by the language model. Results shown in supplementary Fig. 3, Experiment 17.

Summary of experimental approaches in LLM-based note summarisation

Our study evaluated various approaches to prompt design and output structuring for LLM-based clinical notes. The experiments were designed to iteratively refine the model's ability to produce accurate, structured, and clinically relevant summaries. The key methodological approaches across our 18 experiments are summarised below:

Baseline prompts (Experiments 1 and 2). These were the initial prompts we had in our product prior to adopting the framework. We used these as a benchmark against which the later prompts were compared.

Structured prompts (Experiments 3, 7, 8, and 4). Prompts were organised into three components: base (context and goal setting), template (content and structure), and style (formatting). Experiment 3 was a customisation experiment to test a new prompt structure (base, template, style preferences) for generating custom notes based on a transcript. Experiment 7 introduced a first-person perspective in generated notes. Experiment 8 refined the style section by incorporating negations and an "unknown" category for problems not explicitly mentioned in transcripts. Experiment 4 tested an enhancement to the baseline SOAP (Subjective, Objective, Assessment, Plan) note by improving medication record representation.

Atomisation (Experiment 5). This method used a chain-of-thought prompt to extract atomic facts from transcripts to ensure the precise organisation of clinical details. The approach facilitated structured extraction, breaking down information into fundamental components.

Function calls & JSON-based output (Experiments 6, 9, 10, 11). LLMs were instructed to generate responses in structured JSON format instead of free text. This structured format was optimised for integration with primary care electronic health record systems. Successive experiments refined style handling, negation accuracy, and clinical specificity, progressively reducing hallucinations.

Structured prompt + LLM revision step (Experiments 14 and 15). We added a second LLM pass to review and refine outputs based on structured prompting. Experiment 14 built on Experiment 11 with a revision step to improve SOAP notes and introduce an "unknown" option for missing details. Experiment 15 applied this process to a 'Bad SOAP' note, containing hallucinations and omissions, to evaluate how well errors could be mitigated.

New note generation approach (Experiment 16). A novel template-driven method was introduced for generating customised outputs. However, comparison with baseline results (Experiment 8) revealed an increase in major hallucinations and minor omissions, highlighting potential trade-offs.

Clinician vs. LLM comparison (Experiment 17). In this experiment, clinicians manually created notes based on medical transcripts. These notes were then assessed for hallucinations and omissions, providing a benchmark against LLM-generated content. Interestingly, findings suggested slightly more hallucinations in clinician-written notes but fewer omissions, highlighting key differences in human vs. LLM-generated summaries.

Experiment 18. In this experiment, we assessed the performance of the notes within the publicly available ACI Bench dataset.

Hallucination and Omission Taxonomy: We follow the conventional AI literature and taxonomise LLM errors into two types 1) hallucinations, which are instances of text unsupported by the associated clinical documentation, and 2) omissions²¹, which are instances where relevant details are missed in the supporting evidence. Furthermore, inspired by protocols in medical device certifications^{57,58}, we categorise errors as either 'major' or 'minor', where major errors can impact on the diagnosis or the management of the patient if not corrected.

To make our categories more granular, we propose to divide hallucinations into four categories: (1) fabrication, occurring when the model produced information that was not evidenced in the text, (2) negation, occurring when the model output negates a clinically relevant fact, (3) causality, occurring when a model speculates the cause of a given condition without explicit support from the text, and (4) contextual, occurring when the model mixes topics otherwise not related to the given context.

In the case of omissions, we further divide them into sections: (1) current issues, occurring when details about the current presentation were

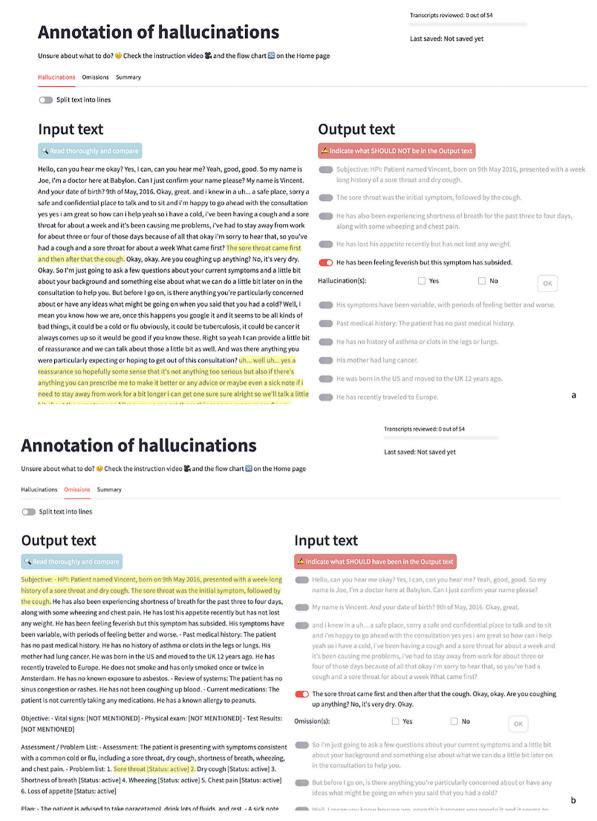


Fig. 15 | The annotation user interface is used to identify hallucinations (a) and omissions (b) and categorise them into major and minor categories. To facilitate clinician review, the closest sentence matches (highlighted in yellow) for each portion of the text under review were extracted from the counterpart document. In the

case of hallucinations, portions of text in the note were compared to the consultation (a), whereas for omissions, portions of the transcript were compared against the note (b).

omitted, (2) PMFS (past medical history, medication history, family and social history), occurring when details about the past medical history, medications including allergies, family and social history, including drinking and smoking, were omitted, and 3) information and plan: when discussions and explanations of the condition and management plans were omitted. Examples of each of the sub-categories are provided in the Supplementary Materials.

Experimental Structure and Annotation Protocol: Here, we define a process to assess how model parameters affect the model outputs and clinical safety. To do this, we define "experiments", which are parametrised by (1) the number of data points processed by the LLM, (2) the type of data the LLM will ingest, (3) the model configuration (type of model, random seed, temperature,...), (4) the prompt used to obtain an LLM output, and (5) the number of clinicians which must review the data point for clinical errors.

Given an experiment configuration, we extract model outputs from the input data and store the results in a database with the associated experiment metadata. We task annotators to classify whether given sub-sections of the output contain hallucinations or omissions according to our taxonomy, and explain in free text the reason for classification. The annotators were volunteer doctors who were paid £5 per note for annotations. Recognising the subjectivity inherent in annotation, we require annotation by at least two clinicians for each input-output pair. This step is followed by a consolidation step, i.e. a detailed review by our internal team of senior clinicians, ensuring a consistent evaluation of all annotations.

Clinical Safety Assessment: Recognising that safety assessment is a crucial part of using any medical technology, we designed a safety evaluation framework of the LLM outputs based on the framework used for evaluating a medical device^{57,58}. Overall, this assessment involves estimating the likelihood of an error happening (Table 3) in conjunction with the potential impact of the error on the clinical outcome if it does occur. Table 4 shows the classification of the level of harm, and Fig. 14 presents the estimation of risk based on the likelihood and consequences of an event.

To maintain consistency in assessing the likelihood of hallucinations and omissions in each experiment, we created a percentage-based metric for their occurrence across experiments, as detailed in Supplementary data 1. 'Very High' likelihood represents scenarios where errors were very common (>90%), whereas 'Very Low' likelihood was associated with situations where errors were rare (<1%). For the 'Medium' likelihood category, which covers 10–60%, we used a broader range to accommodate the output variability and the understanding that some errors may be less predictable or depend on context (Table 5).

CREOLA, Clinical Review of LLMs and AI: We combine the experiment design, hallucination and omission taxonomy, and clinical safety evaluation in a platform we denote CREOLA, short for Clinical Review of LLMs and A1 (pays tribute to Creola Katherine Johnson⁵⁹, a pioneering human computer at NASA. Just as human computers were integral to the safe landing of Apollo moon missions, clinicians play a vital role in safely integrating AI technologies into clinical practice).

The platform is used to identify resultant changes in generated clinical documentation arising from changes to processes in LLM architecture. As illustrated in the "experimental structure", these changes could involve—but are not limited to—the type of model used or prompts used to obtain outputs. The platform was hosted as a Streamlit web application (https://creola.tortus.ai/); the annotation user interface is displayed in Fig. 15.

Annotator recruitment: As outlined earlier, our framework requires annotators to review model outputs. Clinicians are uniquely skilled in critically assessing the veracity of clinical facts in the text. Therefore, we ask clinicians to annotate errors for our experiments. Annotators could register to contribute to the annotation through the CREOLA platform. To ensure annotators had a good understanding of the process, one-to-one tuition was initially provided by the study team. As the number of annotators grew, a short online course was developed to explain the annotation process, followed by a questionnaire to ensure a comprehensive understanding of the material. The annotators were only able to participate if they completed the questionnaire correctly. The annotators could contact the study teams with

any questions through the CREOLA platform in order to ensure any problems in the platform were dealt with promptly.

Date availability

We have used data from Primock and ACI bench which are publicly available clinical transcripts and notes (references in the main text)

Code availability

We have added all our prompts used in the supplementary materials of the article. As explained in the methods section, we used OpenAI's GPT-4 (GPT-4-32k-0613) as the LLM for all our experiments.

Received: 19 August 2024; Accepted: 22 April 2025; Published online: 13 May 2025

References

- Clusmann, J. et al. The future landscape of large language models in medicine. Commun. Med. 3, 141 (2023).
- Becker, G. et al. Four minutes for a patient, twenty seconds for a relative—an observational study at a university hospital. BMC Health Serv. Res. 10, 94 (2010).
- Asgari, E. et al. Impact of electronic health record use on cognitive load and burnout among clinicians: narrative review. *JMIR Med. Inf.* 12, e55499 (2024).
- Ali, S. R., Dobbs, T. D., Hutchings, H. A. & Whitaker, I. S. Using ChatGPT to write patient clinic letters. *Lancet Digit Health* 5, e179–e181 (2023).
- Patel, S. B. & Lam, K. ChatGPT: the future of discharge summaries?. Lancet Digit Health 5, e107–e108 (2023).
- Van Veen, D. et al. Adapted large language models can outperform medical experts in clinical text summarization. *Nat. Med.* 30, 1134–1142 (2024).
- Zhang, T. et al. Benchmarking large language models for news summarization. Trans. Assoc. Comput Linguist 12, 39–57 (2024).
- Ahmad, M. A., Yaramis, I. & Roy, T. D. Creating Trustworthy LLMs: Dealing with Hallucinations in Healthcare Al. Computer Science: Computation and Language. https://doi.org/10.48550/arXiv.2311. 01463 (2023).
- Xu, Z., Jain, S. & Kankanhalli, M. Hallucination is Inevitable: An Innate Limitation of Large Language Models. Computer Science: Computation and Language. https://doi.org/10.48550/arXiv.2401. 11817 (2025).
- Kripalani, S. et al. Deficits in communication and information transfer between hospital-based and primary care physicians. *JAMA* 297, 831 (2007).
- Adane, K., Gizachew, M. & Kendie, S. The role of medical data in efficient patient care delivery: a review. *Risk Manag Health. Policy* 12, 67–73 (2019).
- Schiff, G. D. Diagnostic error in medicine. Arch. Intern Med. 169, 1881 (2009).
- Kumar, S., Balachandran, V., Njoo, L., Anastasopoulos, A. & Tsvetkov, Y. Language generation models can cause harm: so what can we do about it? An actionable survey. In: Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics 3299–3321 (Association for Computational Linguistics, Stroudsburg, PA, USA, 2023). https://doi.org/10.18653/ v1/2023.eacl-main.241.
- Bender, E. M., Gebru, T., McMillan-Major, A. & Shmitchell, S. On the Dangers of Stochastic Parrots. In: Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency 610–623 (ACM, New York, NY, USA, 2021). https://doi.org/10.1145/3442188. 3445922.
- Rando, J. & Tramèr, F. Universal jailbreak backdoors from poisoned human feedback. Computer science: Artificial intelligence. https:// doi.org/10.48550/arXiv.2311.14455 (2024).

- Huang, L. et al. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. Computer science: Computation and language. https://doi.org/10.48550/arXiv. 2311.05232 (2024).
- Rawte, V. et al. Exploring the Relationship between LLM
 Hallucinations and Prompt Linguistic Nuances: Readability,
 Formality, and Concreteness. Computer Science: Artificial
 Intelligence. https://doi.org/10.48550/arXiv.2309.11064 (2023).
- Tonmoy, S. M. T. I. et al. A Comprehensive Survey of Hallucination Mitigation Techniques in Large Language Models. Computer Science, Computation and Language. https://doi.org/10.48550/arXiv. 2401.01313 (2024).
- Overhage, J. M., Qeadan, F., Choi, E. H. E., Vos, D. & Kroth, P. J. Explaining variability in electronic health record effort in primary care ambulatory encounters. *Appl Clin. Inf.* 15, 212–219 (2024).
- Shahbodaghi, A., Moghaddasi, H., Asadi, F. & Hosseini, A. Documentation errors and deficiencies in medical records: a systematic review. *J. Health Manag* 26, 351–368 (2024).
- Moramarco, F. et al. Human Evaluation and Correlation with Automatic Metrics in Consultation Note Generation. In: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) 5739–5754 (Association for Computational Linguistics, Stroudsburg, PA, USA, 2022). https://doi. org/10.18653/v1/2022.acl-long.394.
- Demner-Fushman, D. et al. Preparing a collection of radiology examinations for distribution and retrieval. *J. Am. Med. Inform. Assoc.* 23, 304–310 (2016).
- Abacha, A. Ben, Yim, W., Michalopoulos, G. & Lin, T. An Investigation of Evaluation Metrics for Automated Medical Note Generation. Computer Science: Computation and Language. https://doi.org/10. 48550/arXiv.2305.17364 (2023).
- 24. Ji, Z. et al. Survey of hallucination in natural language generation. *ACM Comput. Surv.* **55**, 1–38 (2023).
- Huang, Y., Tang, K., Chen, M. & Wang, B. A Comprehensive Survey on Evaluating Large Language Model Applications in the Medical Industry. Computer Science: Computation and Language. https://doi. org/10.48550/arXiv.2404.15777 (2024).
- 26. Tierney, A. A. et al. Ambient Artificial Intelligence Scribes to Alleviate the Burden of Clinical Documentation. *NEJM Catal.* **5** (2024).
- Chin-Yew Lin. ROUGE: A Package for Automatic Evaluation of Summaries. In Text Summarization Branches Out, 74–81, Barcelona, Spain. Association for Computational Linguistics. https:// aclanthology.org/W04-1013/ (2004).
- Papineni, K., Roukos. S., Ward. T. & Zhu. W. Bleu: A method for automatic evaluation of machine translation. In: *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, 311–318. Philadelphia, Pennsylvania, USA. https://aclanthology.org/ P02-1040/ (2002).
- Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q. & Artzi, Y. BERTScore: Evaluating Text Generation with BERT. Computer Science: Computation and Language. https://doi.org/10.48550/arXiv. 1904.09675 (2020).
- Moreno, A. C. & Bitterman, D. S. Toward Clinical-Grade Evaluation of Large Language Models. *Int. J. Radiat. Oncol.*Biol.*Phys.* 118, 916–920 (2024).
- Minaee, S. et al. Large Language Models: A Survey. Computer Science: Computation and Language. https://doi.org/10.48550/arXiv. 2402.06196 (2025).
- Evans, O. et al. Truthful Al: Developing and governing Al that does not lie. Computer Science: Computers and Society. https://doi.org/10. 48550/arXiv.2110.06674 (2021).
- Singhal, K. et al. Large Language Models Encode Clinical Knowledge. Computer Science: Computation and Language. https://doi.org/10. 48550/arXiv.2212.13138 (2022).

- Reddy, S. Evaluating large language models for use in healthcare: a framework for translational value assessment. *Inf. Med Unlocked* 41, 101304 (2023).
- 35. Tang, L. et al. Evaluating large language models on medical evidence summarization. *NPJ Digit Med.* **6**, 158 (2023).
- Tam, T. Y. C. et al. A framework for human evaluation of large language models in healthcare derived from literature review. NPJ Digit Med. 7, 258 (2024).
- Cohan, A. et al. A Discourse-Aware Attention Model for Abstractive Summarization of Long Documents (2018).
- Gupta, V., Bharti, P., Nokhiz, P. & Karnick, H. SumPubMed: summarization dataset of pubmed scientific articles. In: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop 292–303 (Association for Computational Linguistics, Stroudsburg, PA, USA, 2021). https://doi.org/10.18653/v1/2021.acl-srw.30.
- Luo, Z., Xie, Q. & Ananiadou, S. CitationSum: Citation-aware Graph Contrastive Learning for Scientific Paper Summarization. *Computer Science: Information Retrieval*. https://doi.org/10.48550/arXiv.2301. 11223 (2023).
- DeYoung, J., Beltagy, I., van Zuylen, M., Kuehl, B. & Wang, L. L. MS2: Multi-Document Summarization of Medical Studies. *Computer Science: Computation and Language*. https://doi.org/10.48550/arXiv. 2104.06486 (2021).
- Song, Y., Tian, Y., Wang, N. & Xia, F. Summarizing Medical Conversations via Identifying Important Utterances. In: Proceedings of the 28th International Conference on Computational Linguistics 717–729 (International Committee on Computational Linguistics, Stroudsburg, PA, USA, 2020). https://doi.org/10.18653/v1/2020.coling-main.63.
- Johnson, A. E. W. et al. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. Sci. Data 6, 317 (2019).
- Pal, A., Umapathi, L. K. & Sankarasubbu, M. Med-HALT: Medical Domain Hallucination Test for Large Language Models. Computer Science: Computation and Language. https://doi.org/10.48550/arXiv. 2307.15343 (2023).
- 44. Korfiatis, A. P. and Moramarco. F. and Sarac. R. and Savkov. A. *PriMock57: A Dataset of Primary Care Mock Consultations*. https:// Github.Com/Babylonhealth/Primock57 (2022).
- Farquhar, S., Kossen, J., Kuhn, L. & Gal, Y. Detecting hallucinations in large language models using semantic entropy. *Nature* 630, 625–630 (2024).
- Simon Hughes, M. B. Vectara Hallucination Leaderboard comparing LLM performance at maintaining factual consistency when summarizing a set of facts. Vectara, Inc. https://huggingface.co/ spaces/vectara/leaderboard (2023).
- 47. Riedemann, L., Labonne, M. & Gilbert, S. The path forward for large language models in medicine is open. *NPJ Digit Med.* **7**, 339 (2024).
- Zhang, G. et al. Closing the gap between open source and commercial large language models for medical evidence summarization. NPJ Digit Med. 7, 239 (2024).
- Lewis, P. et al. Retrieval-augmented generation for knowledgeintensive NLP tasks. Computer Science: Computation and Language. https://doi.org/10.48550/arXiv.2005.11401 (2021).
- Wei, J. et al. Chain-of-thought prompting elicits reasoning in large language models. Computer Science: Computation and Language. https://doi.org/10.48550/arXiv.2201.11903 (2023).
- Jia, M., Duan, J., Song, Y. & Wang, J. medlKAL: integrating knowledge graphs as assistants of LLMs for enhanced clinical diagnosis on EMRs. Computer Science: Computation and Language. https://doi. org/10.48550/arXiv.2406.14326 (2025).
- Wang, Y., Ma, X. & Chen, W. Augmenting black-box LLMs with medical textbooks for clinical question answering. Computer

- Science: Computation and Language. https://doi.org/10.48550/arXiv. 2309.02233 (2024).
- Gilbert, S., Kather, J. N. & Hogan, A. Augmented non-hallucinating large language models as medical information curators. NPJ Digit Med. 7, 100 (2024).
- Desmond, M., Ashktorab, Z., Pan, Q., Dugan, C. & Johnson, J. M. EvaluLLM: LLM assisted evaluation of generative outputs. In: Companion Proceedings of the 29th International Conference on Intelligent User Interfaces 30–32 (ACM, New York, NY, USA, 2024). https://doi.org/10.1145/3640544.3645216.
- Zheng, L. et al. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. Computer Science: Computation and Language. https://doi. org/10.48550/arXiv.2306.05685 (2023).
- Gu, J. et al. A survey on LLM-as-a-Judge. Computer Science: Computation and Language. https://doi.org/10.48550/arXiv.2411. 15594 (2024).
- International Organization for Standardization. https://Www.lso.Org/ Standard/59752.Html.
- International Organization for Standardization. https://Www.lso.Org/ Standard/72704.Html.
- https://en.wikipedia.org/wiki/Katherine_Johnson. Katherine Johnson.

Acknowledgements

The authors would like to thank the many physicians who helped with our CREOLA experiments and the annotation of the clinical transcripts and notes. No funding has been received for this study.

Author contributions

D.P., E.A., M.D., N.M., S.K. and J.B. contributed to the concept, design and execution of the study. M.D. and S.K. built the CREOLA platform and with NM designed the various prompts and experiments. MD analysed all the results and prepared the figures. E.A. and N.M. wrote, reviewed and revised the paper. E.A. and D.P. designed the clinical safety framework, reviewed all the annotations and scored the impact of errors on clinical safety. JAY provided expert advice on paper writing and direction and contributed to significant paper revision and rewriting. E.A., J.A.Y., M.D., N.M., S.K., J.B. and D.P. contributed to the review and revision of the paper.

Competing interests

D.M. is the CEO of the company Tortus AI and all authors were employees of Tortus AI at the time of the writing the paper (M.D., N.M., S.K. and J.B. as full time and E.A. part time, J.A.Y. full time employee at the time of revision). The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at https://doi.org/10.1038/s41746-025-01670-7.

Correspondence and requests for materials should be addressed to Elham Asgari.

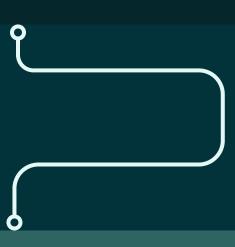
Reprints and permissions information is available at http://www.nature.com/reprints

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit http://creativecommons.org/licenses/by-nc-nd/4.0/.

© The Author(s) 2025

Whether you're exploring ambient AI, planning a pilot, or shaping digital strategy, we'd love to talk.

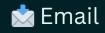




Book a meeting with our CEO and founder



Dr. Dom Pimenta



dom@tortus.ai

