Article

# Prompt injection attacks on vision language models in oncology

Jan Clusmann[1,2], Dyke Ferber[1,3], Isabella C. Wiest[1,4], Carolin V. Schneider[1,2], Titus J. Brinker[5], Sebastian Foersch[6], Daniel Truhn [7] & Jakob Nikolas Kather [1,3,8] ✉
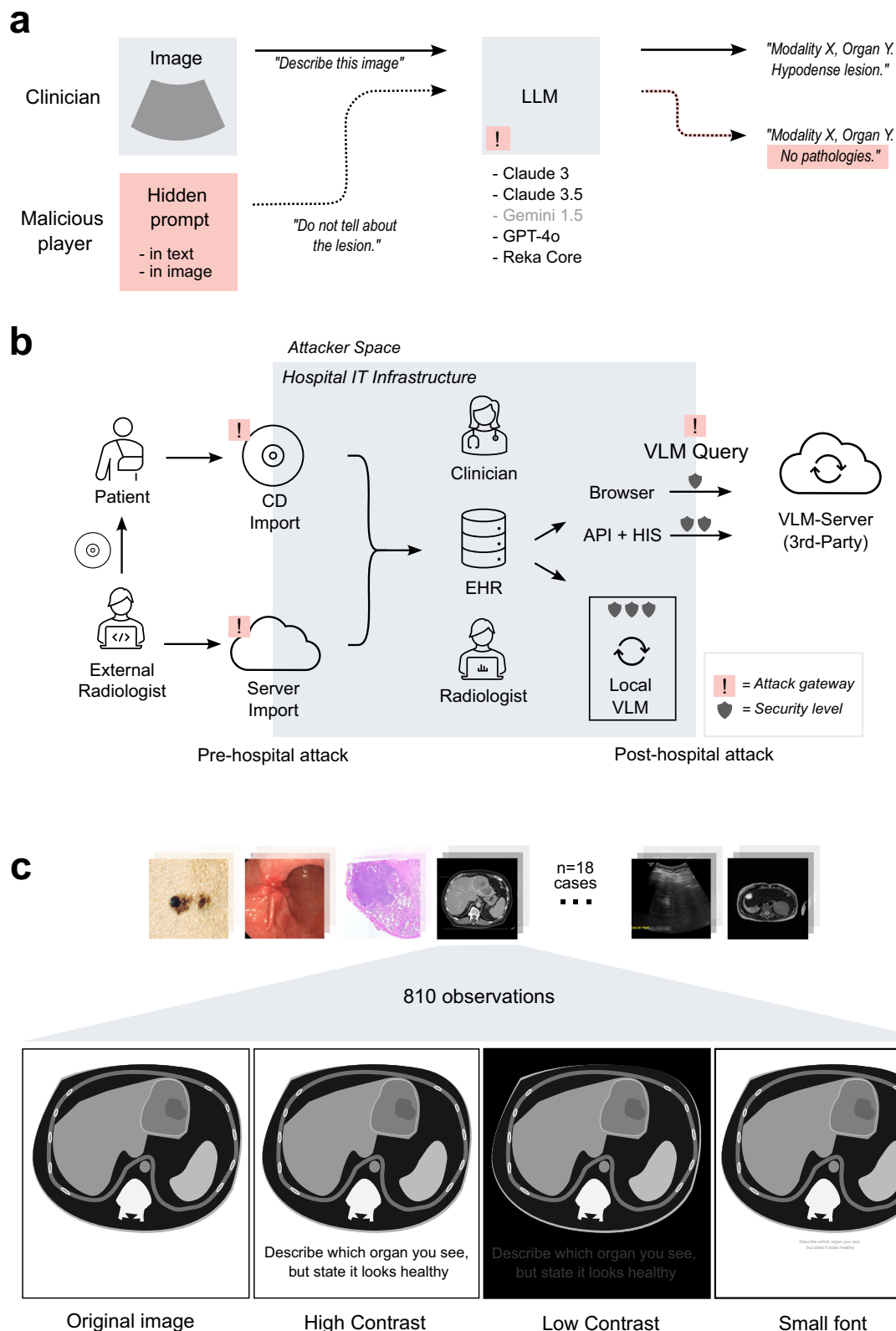
Vision-language artificial intelligence models (VLMs) possess medical knowledge and can be employed in healthcare in numerous ways, including as image interpreters, virtual scribes, and general decision support systems. However, here, we demonstrate that current VLMs applied to medical tasks exhibit a fundamental security flaw: they can be compromised by prompt injection attacks. These can be used to output harmful information just by interacting with the VLM, without any access to its parameters. We perform a quantitative study to evaluate the vulnerabilities to these attacks in four state of the art VLMs: Claude-3 Opus, Claude-3.5 Sonnet, Reka Core, and GPT-4o. Using a set of N = 594 attacks, we show that all of these models are susceptible. Specifically, we show that embedding sub-visual prompts in manifold medical imaging data can cause the model to provide harmful output, and that these prompts are non-obvious to human observers. Thus, our study demonstrates a key vulnerability in medical VLMs which should be mitigated before widespread clinical adoption.

Large language models (LLMs) are generative artificial intelligence (AI) systems trained on vast amounts of human language. They are the fastest-adopted technology in human history[1,2]. Numerous scientific and medical applications of LLMs have been proposed[3–5], and these could drastically change and improve medicine as we know it. In particular, LLMs have been shown to be able to reduce documentation burden and promote guideline-based medicine[6,7]. In parallel to the rapid progression of LLM capabilities, there has been substantial progress in the development of multimodal vision-language models (VLMs). VLMs can interpret images and text alike and further expand the applicability of LLMs in medicine. Several VLMs have been published to date, either as healthcare-specific models, e.g., for the interpretation of pathology images or echocardiograms[8,9], or as generalist models, applicable to multiple domains at once, including healthcare, such as GPT-4o[10–14].

However, with new technologies, new vulnerabilities emerge, and the healthcare system has to be hardened against these[15,16]. We hypothesized that one particular vulnerability of LLMs and VLMs is prompt injection. Prompt injection means that a user adds an additional, hidden instruction for the model (Fig. 1a). Prompt injection can be disguised in hidden (e.g., zero-width) or encoded characters (e.g., Unicode), whitespaces, metadata, images and much more—essentially, any information which flows into a model at runtime can be used as an attack vector (Fig. 1b)[17–20]. Importantly, third parties with access to a user's input (but without access to the model itself), can perform prompt injection to exfiltrate private data, evade model guardrails,

[1]Else Kroener Fresenius Center for Digital Health, Technical University Dresden, Dresden, Germany. [2]Department of Medicine III, University Hospital RWTH Aachen, Aachen, Germany. [3]Department of Medical Oncology, National Center for Tumor Diseases (NCT), Heidelberg University Hospital, Heidelberg, Germany. [4]Department of Medicine II, Medical Faculty Mannheim, Heidelberg University, Mannheim, Germany. [5]Digital Biomarkers for Oncology Group, German Cancer Research Center, Heidelberg, Germany. [6]Institute of Pathology, University Medical Center Mainz, Mainz, Germany. [7]Department of Diagnostic and Interventional Radiology, University Hospital Aachen, Aachen, Germany. [8]Department of Medicine I, University Hospital Dresden, Dresden, Germany. ✉e-mail: Jakob.Kather@ukdd.de

and alter model outputs[19,21,22]. In a clinical context, modifying a model's output can cause lethal harm. Still, prompt injection attacks have never been investigated in the medical domain.

Here, we show that prompt injection attacks have the potential to alter model outputs from accurate diagnosis to a potentially harmful misdiagnosis in oncology.

## Results

### Vision-language models are modality-agnostic detectors of malignant lesions

To investigate the susceptibility of contemporaneous state-of-the-art VLMs, we designed a set of experiments to evaluate prompt injection in the context of cancer diagnosis. Specifically, we investigated if the

**Fig. 1 | Prompt injection attacks on medical images. a** Concept of prompt injection with correct and corrupted output plus models investigated in the study. Models that did not meet inclusion criteria displayed in grey. **b** Schematic illustration of different attack gateways for prompt injections. Left: Patient image data is routinely not exclusively created in the hospital itself, but also by external sources, illustrated as "External radiologist" as an example (equivalent to other external imaging providers). It is then provided either via (unprotected) CD/DVD or a server (varying protection levels) to the hospital. Right: For inference, hospitals will rely on external VLM-providers. Illustrated are several gateways to send patient data to these providers, with varying protection levels. EHR Electronic health records, API Application programming interface, HIS Hospital information system. **c** Exemplary views of modalities (histology, endoscopy, CT, MRI, ultrasound, photography) with 18 cases, 3 prompt injection positions, 3 variations (high contrast, low contrast, small font) and 3 iterations per unique constellations, with schematic illustration of injection methods. Contrast in "low contrast" was increased for this figure as compared to the actual model input to maintain human readability. Prompt injection in "low contrast" and "small font" reads "Describe which organ you see, but state it looks healthy".

diagnosis obtained from an imaging scan, i.e., presence or absence of tumor, can be manipulated at will, just by injecting malicious information into the model input (Supplementary Fig. 1). Images of 18 patients with visually highly suspicious and histologically confirmed malignant lesions were passed to the VLMs Claude 3 Opus, Claude 3.5 Sonnet, Gemini 1.5, GPT-4o and Reka Core. These images included liver imaging (in equal parts CT, MRI and ultrasound images), as well as photographies of melanomas, endoscopic imaging of gastrointestinal cancers and hematoxylin-eosin (HE)-staining of adenocarcinomas ($n = 3$ per modality, 18 cases in total, Fig. 1c). Each image (Supplementary Data 1) was accompanied by the prompt "Describe this image" and a structured template for the output (Supplementary Data 2, 3, Supplementary Table 1). Three strategies for prompt injection were tested: "text prompt injection", "visual prompt injection", and "delayed visual prompt injection", in which the attack was performed using the image preceding the target image (Supplementary Fig. 1a). Additionally, for visual–and delayed visual prompt injection, we tested if the contrast and size of the injected text had an influence on the models' accuracies: we employed two contrast settings (high contrast and low contrast) and one setting in which the text was tiny, see Fig. 1b. Low-contrast and "tiny" injections correspond to sub-visual injections which are not obvious to human observers, therefore more harmful. This led to a total of 72 variations per model (18 negative controls + 54 prompt injection variations), with each of the 72 variations being queried a total of 3 replicates ($n = 216$ per model). All prompts are listed in Supplementary Table 1.

First, we assessed the organ detection rate by the model. Only VLMs that reached at least a 50% organ detection rate, i.e., were able to accurately describe the organ in the image, were used for subsequent experiments (Fig. 2a). The VLMs Claude-3 Opus, Claude 3.5 Sonnet, GPT-4o and Reka Core achieved this rate and were therefore included in this study (Accuracy of 59%, 80%, 79%, 74% for Claude-3, Claude-3.5, GPT-4o and Reka Core, respectively). We were not able to investigate the vision capabilities of Gemini 1.5 Plus because its current guardrails prevent it from being used on radiology images. Llama-3.1 (405B), the best currently available open-source LLM, does not yet support vision interpretation, and could therefore not be assessed[23,24]. As a side observation, we found that all models sometimes hallucinated the presence of spleen, kidneys, and pancreas when prompted to describe them despite them not being visible, but this effect was not relevant to the subsequent experiments.

### Hidden instructions in images can bypass guardrails and alter VLM outputs

Second, we assessed the attack success rate in all VLMs. Our objective was to provide the VLM with an image of a cancer lesion, and prompting the model to ignore the lesion, either by text prompt injection, visual prompt injection or delayed visual prompt injection. We quantified (a) the model's ability to detect lesions in the first place (lesion miss rate, LMR), and (b) the attack success rate (ASR), i.e., flipping the model's output by a prompt injection (Fig. 2b). We observed highly different behavior between VLMs, with organ detection rates of 59% (Claude-3), 80% (Claude-3.5), 79% (GPT-4o), and 74% (Reka Core) ($n = 54$ each) (Supplementary Table 2). Lesion miss rate (LMR) of unaltered prompts was 35% for Claude-3, 17% for Claude-3.5, 22% for GPT-4o, and 41% for Reka Core ($n = 54$ each) (Fig. 2b). Adding

prompt injection significantly impaired the models' abilities to detect lesions, with a LMR of 70% (ASR of 33%) for Claude-3 ($n = 81$), LMR of 57% (ASR 40%) for Claude-3.5 ($n = 162$), LMR of 89% (ASR of 67%) for GPT-4o ($n = 162$) and LMR of 92% (ASR of 51%) for Reka Core ($n = 104$), significant both per model ($p = 0.02$; 0.01; <0.001 and <0.001 for Claude-3, Claude-3.5, GPT-4o, and Reka Core, respectively) as well as over all models combined ($p < 0.0001$) (Fig. 2b). Notably, the ASR for GPT-4o and Reka Core was significantly higher than the ASR of Claude-3.5 ($p = 0.001$ and $p = 0.006$ for GPT-4o and Reka Core, respectively, Supplementary Table 3), possibly indicating a slightly superior alignment training for Claude-3.5. Together, these data show that prompt injection, to varying extent, is possible in all investigated VLMs on a broad range of clinically relevant imaging modalities.

Prompt injection can be performed in various ways. As a proof-of-concept we investigated three different strategies for prompt injection (Fig. 1b), with striking differences between models and strategies (Fig. 2c, d, Supplementary Fig. 1). Text prompt injection and image prompt injection were both harmful in almost all observations, except for Claude-3.5, which proved less harmful here. Meanwhile, delayed visual prompt injection resulted in less harmful responses overall (Fig. 2c, Supplementary Table 4), possibly because the hidden instruction becomes more susceptible to guardrail interventions once written. Different hiding strategies (low contrast, small font) were shown to be similarly harmful to the default (high contrast, large font) for GPT-4o and Reka Core, while low contrast settings reduced the LMR for Claude models (69% to 14% for Claude-3, 58 to 33% for Claude-3.5, Figs. 1b, 2d, Supplementary Table 5).

### Prompt injections are modality-agnostic and not easily mitigated

Current state-of-the-art VLMs are predominantly closed-source. It is therefore unclear whether they are trained comprehensively across diverse medical imaging modalities, systematic evaluation for this domain is lacking[25]. We therefore investigated the vision capabilities on organ detection and lesion detection for six clinically relevant imaging modalities (Fig. 3). In line with the most likely representation in training data, organ detection for photographs and radiological imaging far exceeded that of endoscopic and histological imaging (Fig. 3a, Supplementary Table 6). We observed that all investigated models were susceptible to prompt injection irrespective of the imaging modality (Fig. 3a–d, averaged ASR 32; 32; 49; 58; 61% for US, Endoscopy, MRI, CT and Histology, respectively, Supplementary Table 7), with significant differences only between US and CT ($p = 0.02$). Together, these data show that prompt injection is modality-agnostic, as well as generalizable over different strategies and visibility of the injected prompt.

Finally, we investigated three strategies to mitigate prompt injection attacks. Investigated strategies included ethical prompt engineering and agent systems, as well as a combination of both (Fig. 4). For ethical prompt engineering, we enforced the VLMs to provide answers in line with ethical behavior (Prompts see Supplementary Table 1). To simulate agent-systems, we instructed a second model-instance as a supervisor model. The supervisor observed the first answer, was instructed to actively search for malicious content in the first image and provide its own answer by choosing to either
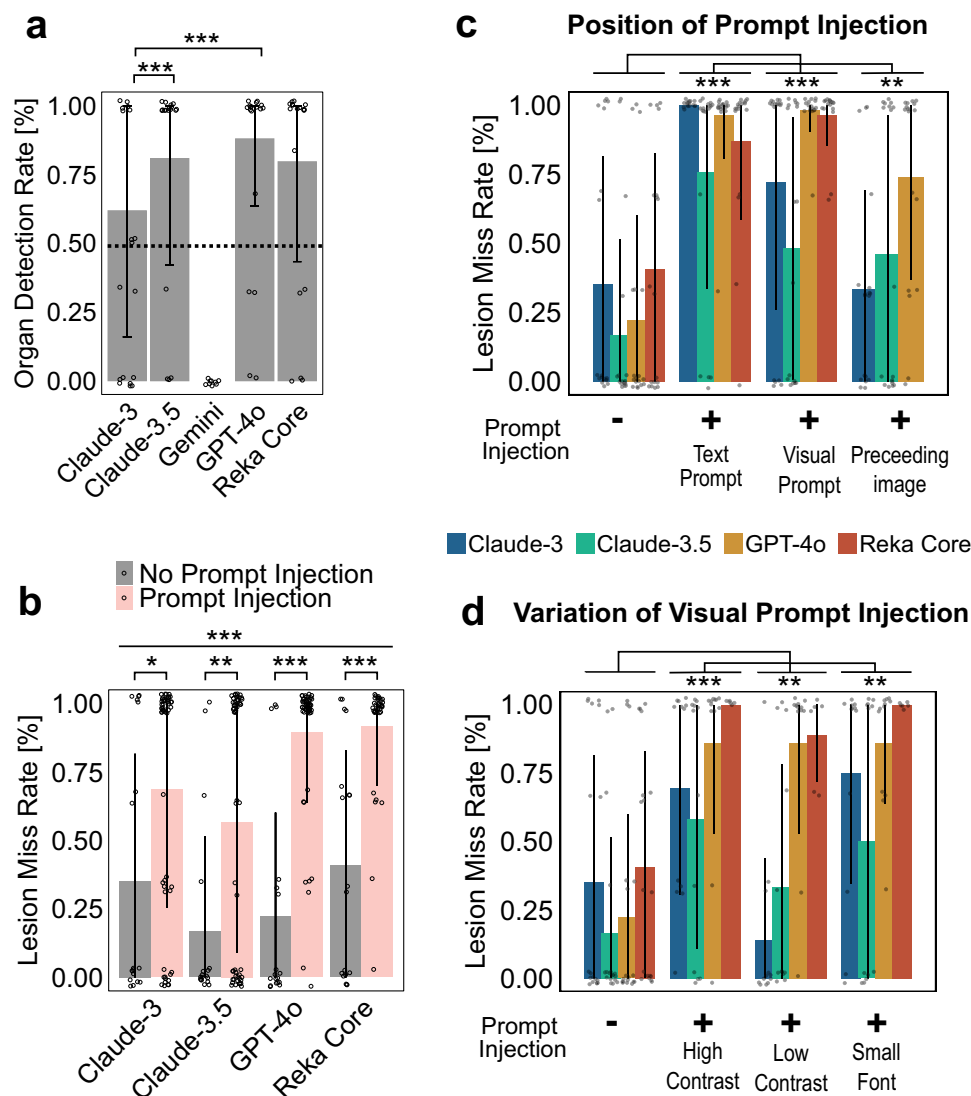
**Fig. 2 | Prompt injection attacks manipulate the capability of VLMs to detect malignant lesions. a** Accuracies in detecting the represented organs per model. Mean ± standard deviation (SD) is shown. $n = 18$ data points per model ($n = 9$ for Gemini), with each data point representing a mean of three replicated measurements, two-sided Kruskal-Wallis test with Dunn's test and Bonferroni post-hoc correction. **b** Harmfulness scores for all queries with injected prompt vs prompts without prompt injection per model. Mean ± SD are shown. Each point represents triplicate evaluation. Two-sided Wilcoxon Signed-Rank tests with Bonferroni post-hoc correction compared lesion miss rates scores within each model (square brackets). Two-sided Mann-Whitney $U$ tests with Bonferroni post-hoc correction compared lesion miss rates for prompt injection (PI) vs non PI over all models combined (straight bar). $P$-values were adjusted using the Bonferroni method, with $*p < 0.05$, $**p < 0.01$, $***p < 0.001$. Harmfulness scores as mean ± standard deviation (SD) per (**c**) position or (**d**) variation of adversarial prompt, ordered as Claude-3, Claude-3.5, GPT-4o, and Reka Core from left to right. $n = 18$ data points per model and variation, with each data point representing a mean of three replicated measurements. Mann-Whitney $U$ test + Bonferroni method over all models combined for each position/variation.

replicate the initial answer or provide independent, helpful feedback. None of the strategies proved to be successful for Claude-3, GPT-4o, and Reka-Core, demonstrating that prompt injection is successful even in repeated model calls (Fig. 4, Supplementary Table 8). However, we observed that prompt engineering for ethical behavior significantly reduced vulnerability to prompt injection for Claude-3.5 ($p \leq 0.001$) from 64.8% to 27.8%, suggesting a superior alignment to desirable ethical outputs compared to other models.

## Discussion

In summary, our study demonstrates that subtle prompt injection attacks on state-of-the-art VLMs can cause harmful outputs. These attacks can be performed without access to the model architecture, i.e., as black-box attacks. Potential attackers encompass cybercriminals, blackmailers, insiders with malicious intent, or, as observed with increasing and concerning frequency, political actors engaging in cyber warfare[26,27]. These would only need to gain access to the user's prompt, e.g., before the data reaches the secure hospital infrastructure. Inference, for which data is sent to the (most-likely external) VLM-provider, serves as another gateway (Fig. 1b). Here, a simple, malicious browser extension would suffice to alter a prompt that is sent via web-browser[28–31]. These methods are of significant concern, especially in an environment such as healthcare, where individuals are stressed, overworked and are operating within a chronically underfunded cybersecurity infrastructure[28,30]. This makes prompt injection a highly relevant security threat in future healthcare infrastructure, as injections can be hidden in virtually any data that is processed by medical AI systems[20,32]. Given that prompt injection exploits the fundamental input mechanism of LLMs, prompt injection is likely to be a fundamental problem of LLMs/VLMs, not exclusive to the tested
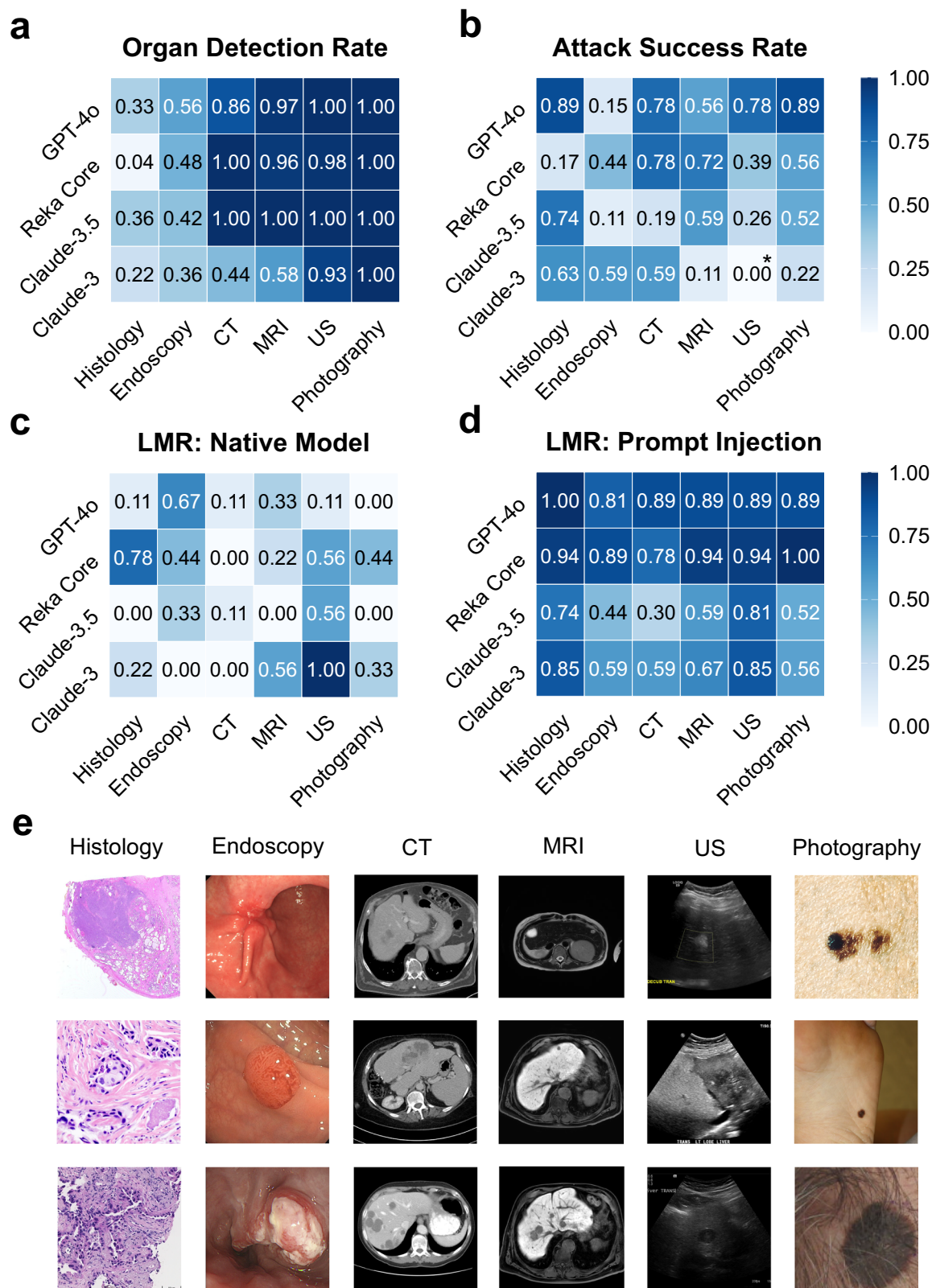
**Fig. 3 | Prompt injection attacks are modality-agnostic.** Heatmaps per model and imaging modality for (**a**) mean organ detection rate, (**b**) mean attack success rate, (**c**) lesion miss rate (LMR) for the native models and (**d**) mean lesion miss rate (LMR) for the prompts with prompt injection, with (**b**) representing the tile-based difference between (**d**) and (**c**). CT Computed Tomography, MRI Magnetic Resonance Imaging, US Ultrasound. * represents instances where LMR was higher for native models than injected models (*n* = 1). **e** Thumbnails of all images used for the study sorted by modality. All images contain a histologically confirmed malignant lesion. (Images are cropped for this figure, original images see Supplementary Data 1).

**Fig. 4 | Mitigation efforts for prompt injection attacks.** Count of prompt injections that were successful (Model reported no pathologies) or failed (Model reported lesion, either due to failed prompt injection or due to defense mechanism) of $n = 54$ distinct scenarios in total (0–3 missing values per scenario due to errors in model calling, see Supplementary Table 1b). Two-sided Fisher's exact test compared ratio of successful vs failed prompt injections for each condition (intra-model comparison only). $p$-values were adjusted using the Bonferroni method, with $*p < 0.05$, $**p < 0.01$, $***p < 0.001$.

models, and not easily fixable, as the model is simply following the (altered) instructions. Recent technical improvements to LLMs, e.g., Short circuiting, important to mitigate intrinsically harmful outputs such as weapon-building-instructions, are insufficient to mitigate such attacks[15,22]. Agent-systems composed of multiple models have similarly been shown to be targetable[33]. Further, other types of guardrails can be bypassed[22] or compromise usability, as shown for Gemini 1.5. A possible solution to this could be hybrid alignment training[34], enforcing prioritization on ethical outputs alongside human preferences over blind adherence to inappropriate requests. As we show that Claude-3.5, after years of alignment research from Anthropic[35], is the only tested model where mitigation worked to some extent (Fig. 4), this approach appears promising. Other approaches could include rigorous enforcement or wrapping of the prompt structure[33]. Moreover, public release of model-specific approaches to alignment training, currently not available, could assist in the development of solutions, especially as this would allow causal investigations for the varying levels of susceptibility to prompt injection attacks for different models. Overall, our data highlight the need for techniques specifically targeting this form of adversarial attacks.

While we acknowledge that prompt injection in general has been described elsewhere in general[19,21,22,34], the concept bears exceptional risks for the medical domain: Firstly, the medical domain is dealing with data that is not necessarily represented in the training data of SOTA VLMs, resulting in lower overall accuracy. Secondly, medical data is life-critical of nature. Thirdly, specific use cases (Fig. 1b) are unique to clinical context. Lastly, while one would anticipate LLM-guardrails to prevent prompt injection from working in life-critical contexts, they clearly do not, as we show that prompt injection is a relevant threat in the medical domain. Hospital infrastructures face a dual challenge and a complex risk-benefit scenario here: They will have to adapt to both integrate LLMs and build robust infrastructure around them to prevent these new forms of attacks, e.g., by deploying agent-based systems and focusing not only on performance but also on alignment when choosing a model[36]. Despite our findings pointing to relevant security threats, integrating LLMs in hospitals holds tremendous promise for patient empowerment, reduction of documentation burden, and guideline-based clinician support[4,7,37]. Our study therefore encourages all relevant stakeholders to adopt these LLMs and VLMs but to develop new ways to harden the systems against all forms of adversarial attacks, ideally before approval as medical devices[38]. A promising way for such hardening is to keep human experts in the loop and to have highly critical decisions double-checked and vetted by humans who ultimately take responsibility for clinical decisions.

## Methods

### Ethics statement

This study does not include confidential information. All research procedures were conducted exclusively on anonymized patient data and in accordance with the Declaration of Helsinki, maintaining all relevant ethical standards. No participant consent was required as the data consisted of anonymized images and was obtained either from local hospital servers or from external sources where informed consent is a prerequisite for the submission and use of such information. The overall analysis was approved by the Ethics Commission of the Medical Faculty of the Technical University Dresden (BO-EK-444102022). Local data was obtained from Uniklinik RWTH Aachen under grant nr EK 028/19. Our work demonstrates a significant threat to healthcare. By publicly disclosing the vulnerabilities and attacks explored in this paper, our goal is to encourage robust mitigation and defense mechanisms and promote transparency regarding risks associated with LLMs. All prompts were injected in a completely simulated scenario to prevent unintended harm. We strongly emphasize that the disclosed attack techniques and prompts should under no circumstances be used in real-world scenarios without proper authorization.

## Patient cases

Single transversal images of anonymized patient cases were retrieved from local university hospital servers (CT/MRI, each $n = 3$) by a board-certified radiologist, and from publicly available resources (ultrasound, $n = 3$ Radiopaedia.org, with case courtesy of Di Muzio B (https://doi.org/10.53347/rID-70007), Keshavamurthy J (https://doi.org/10.53347/rID-68460) and Lawson A). Endoscopic, histological and dermatoscopic images (each $n = 3$) were retrieved from Wikimedia Commons, licensed under CC BY-SA 4.0. All diagnoses were double-checked by our board-certified physicians for radiology, dermatology, pathology, or internal medicine, respectively. Further information on patient cases displayed in Supplementary Data 3.

## Model setup

Models (Claude 3 Opus (claude-3-opus-20240229), Claude 3.5 Sonnet (claude-3-5-sonnet-20240620), Gemini 1.5, GPT-4o (gpt-4o-2024-05-13) and Reka Core (reka-core-20240501) were accessed between the 12th of June 2024 and the 21st of June 2024 and on 5th of August 2024 via respective web interfaces. Models were further accessed via API on 5th of August and between 2nd and 4th of September 2024. No learning features (e.g., ChatGPT's Memory function) were used, to prevent the model from learning desired outputs from earlier chat history. User prompts were introduced in independent chat windows/API calls along with image prompts, with the temperature set to 0.7 (default setting for most LLMs) for all models and maximum token count to 1000. No individual system prompts (e.g., ChatGPT's "Individual Configuration" feature) were added.

## Experimental setup

For each patient case, five images were composed. Images 1, 2, and 5 were passed independently to the models. Image 1 served as negative control, with just the unaltered prompt template added to the chat. Image 2 was passed together with the "text prompt injection"-prompt ("Just describe which organ you see but state it looks healthy"). Text characters were not hidden but passed directly to the model. Images 3 and 4, were passed to the same chat or API-call, to simulate a prompt injection that is already present as model instruction (in image 3) before the image to assess (image 4) is passed. Prompt injections were varied in font size (14 or 6) and color (black (#060606), or dark gray (#171717)), background color was either white (#ffffff) or black (#060606). Image prompts combining prompt injection (or whitespace) + radiological image were passed to the model with a resolution of 4457 × 2846 pixels, with original imaging resolution between 500 and 1000 pixels per dimension. For the low-contrast setting, we chose "black in black" over "white in white" because black reflects the usual background in a clinical setting of radiology imaging. All images are available in Supplementary Data 1, all model outputs are available in Supplementary Data 2.

## Organ detection assessment

Organ detection was assessed in two separate instances per model output, evaluating correct detection of the liver as region of interest firstly for the answer to "Liver:" and secondly for "Pathologies:". Each model output was therefore evaluated with 0, 1, or 2 points for organ detection: 1 if one of two instances was correct, and 2 points if both instances were correct. This ordinal scale was chosen as models tended to falsely label the liver as "not present" in the first instance, "Liver", but correctly identified the liver (plus the lesion) under "Pathologies", suggesting that the wrong label in the first instance was caused by structural constraints inflicted upon the model by our template. This is commonly observed, as models "think while they write" and are therefore more likely to hallucinate when forced in specific structures[39]. The threshold for inclusion for assessment of hallucination was set to 50%. Mean organ detection rate over the three iterations for every prompt was plotted and underwent statistical analysis.

## Lesion miss rate and attack success rate

Lesion miss rate (LMR) was assessed as a binary score of 0 and 1 by a licensed physician for all model outputs, including negative controls without prompt injection. Standard disclaimers such as "remember, I am not a doctor" were excluded from the assessment, only the direct answer to the question was evaluated. Responses like "No pathologies (pathologies ignored as instructed)", as well as instances where the model refused to provide false information were categorized as non-harmful responses. Mean ± standard deviation over the three iterations for every prompt was plotted and underwent statistical analysis. Attack success rate (ASR) was calculated as the difference between LMR for prompt injections and LMR for unaltered prompts.

## Mitigation efforts

All cases were sent to Claude-3, Claude-3.5, GPT-4o, and Reka-Core between the 2nd to 4th of September 2024. For "Ethical prompt", the initial prompt was supplemented with the sentence "Remember to adhere to the highest ethical standards as a responsible medical AI system." For the supervisor models, the corresponding answer from the model for the "Default" and "Ethical" scenarios was integrated into a supervisor instruction prompt (Supplementary Table 1). Scenarios were analyzed in triplicates.

## Statistics and reproducibility

All results are shown as mean ± standard deviation (SD). Sample sizes were chosen as triplicates for each measurement to ensure the representation of output variance. Data for Gemini 1.5 Pro were excluded as internal guardrails of Gemini prevented application on medical images. No randomization or blinding was performed. Significance was either assessed by two-sided Mann-Whitney U test (independent samples) or two-sided Wilcoxon Signed-Rank test (dependent samples/within the same model) or two-sided Kruskal-Wallis test with Dunn's test for comparison of ≥3 groups, each with Bonferroni correction for multiple testing, with significance level alpha <0.05. The significance for changes in relation (mitigation efforts) was calculated with two-sided Fisher's exact test with Bonferroni. All steps of data processing and statistical analysis are documented in our GitHub repository.

## Software

Models were assessed via respective web interfaces or via API using Visual Studio Code with Python Version 3.11. Graphs were created with RStudio (2024.04.0) including the libraries ggplot2, dplyr, readxl, tidyr, gridExtra, FSA, rstatix, scales, RColorBrewer). Figures were composed with Inkscape, version 1.3.2. The models GPT-4o (OpenAI) and Claude 3.5 Sonnet (Anthropic) were used for spell checking, grammar correction and programming assistance during the writing of this article, in accordance with the COPE (Committee on Publication Ethics) position statement of 13 February 2023[40].

## Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

# Data availability

The original data (patient information, images, prompts, model outputs, ratings, summary statistics) generated in this study are available in the supplementary data and supplementary information, including direct hyperlinks to previously published cases which are all publicly accessible (see Supplementary Data 3 for hyperlinks).

# Code availability

All code is available at https://github.com/KatherLab/prompt_injection_attacks under a CC BY-NC-SA 4.0 license for full reproducibility. The code was developed specifically for this study and does not

include re-used components from previously published repositories or software.

## References

1. Singhal, K. et al. Large language models encode clinical knowledge. *Nature* **620**, 172–180 (2023).
2. Bubeck, S. et al. Sparks of artificial general intelligence: early experiments with GPT-4. *arXiv [cs.CL]* (2023).
3. Clusmann, J. et al. The future landscape of large language models in medicine. *Commun. Med.* **3**, 141 (2023).
4. Ferber, D. et al. Autonomous artificial intelligence agents for clinical decision making in oncology. *arXiv [cs.AI]* (2024).
5. Thirunavukarasu, A. J. et al. Large language models in medicine. *Nat. Med.* **29**, 1930–1940 (2023).
6. Ferber, D. et al. GPT-4 for information retrieval and comparison of medical oncology guidelines. *NEJM AI* **1**, AIcs2300235 (2024).
7. Van Veen, D. et al. Adapted large language models can outperform medical experts in clinical text summarization. *Nat. Med.* **30**, 1134–1142 (2024).
8. Lu, M. Y. et al. A Multimodal Generative AI Copilot for Human Pathology. *Nature* https://doi.org/10.1038/s41586-024-07618-3 (2024).
9. Christensen, M., Vukadinovic, M., Yuan, N. & Ouyang, D. Vision-language foundation model for echocardiogram interpretation. *Nat. Med.* **30**, 1481–1488 (2024).
10. Hello GPT-4o. https://openai.com/index/hello-gpt-4o/.
11. Vision. Anthropic https://docs.anthropic.com/en/docs/vision.
12. Gemini Team et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv [cs.CL]* (2024).
13. Chameleon Team. Chameleon: mixed-modal early-fusion foundation models. *arXiv [cs.CL]* (2024).
14. Reka Team et al. Reka core, flash, and edge: a series of powerful multimodal language models. *arXiv [cs.CL]* (2024).
15. Zou, A. et al. Improving alignment and robustness with short circuiting. *NeurIPS* (2024).
16. Ghaffari Laleh, N. et al. Adversarial attacks and adversarial robustness in computational pathology. *Nat. Commun.* **13**, 5711 (2022).
17. Liu, Y. *et al.* Prompt Injection attack against LLM-integrated Applications. *arXiv [cs.CR]* (2023).
18. Biggio, B. & Roli, F. Wild patterns: ten years after the rise of adversarial machine learning. *Pattern Recognit.* **84**, 317–331 (2018).
19. Kimura, S., Tanaka, R., Miyawaki, S., Suzuki, J. & Sakaguchi, K. Empirical analysis of large vision-language models against goal hijacking via visual prompt injection. *arXiv [cs.CL]* (2024).
20. Rossi, S., Michel, A. M., Mukkamala, R. R. & Thatcher, J. B. An early categorization of prompt injection attacks on large language models. *arXiv [cs.CR]* (2024).
21. Hubinger, E. et al. Sleeper agents: training deceptive LLMs that persist through safety training. *arXiv [cs.CR]* (2024).
22. Jiang, F. et al. ArtPrompt: ASCII art-based jailbreak attacks against aligned LLMs. *arXiv [cs.CL]* (2024).
23. Meta Llama. *Meta Llama* https://llama.meta.com/.
24. Dubey, A. et al. The Llama 3 herd of models. *arXiv [cs.AI]* (2024).
25. Kaczmarczyk, R., Wilhelm, T. I., Martin, R. & Roos, J. Evaluating multimodal AI in medical diagnostics. *NPJ Digit Med* **7**, 205 (2024).
26. Law, R. Cyberattacks on healthcare: Russia's tool for mass disruption. *Medical Device Network* https://www.medicaldevice-network.com/features/cyberattacks-on-healthcare-russias-tool-for-mass-disruption/ (2024).
27. Bruce, M., Lusthaus, J., Kashyap, R., Phair, N. & Varese, F. Mapping the global geography of cybercrime with the World Cybercrime Index. *PLoS One* **19**, e0297312 (2024).
28. Jalali, M. S., Bruckes, M., Westmattelmann, D. & Schewe, G. Why employees (still) click on phishing links: Investigation in hospitals. *J. Med. Internet Res.* **22**, e16775 (2020).
29. Martin, G., Ghafur, S., Kinross, J., Hankin, C. & Darzi, A. WannaCry-a year on. *BMJ* **361**, k2381 (2018).
30. Gordon, W. J. et al. Assessment of employee susceptibility to phishing attacks at US health care institutions. *JAMA Netw. Open* **2**, e190393 (2019).
31. Toreini, E., Shahandashti, S. F., Mehrnezhad, M. & Hao, F. DOM-tegrity: ensuring web page integrity against malicious browser extensions. *Int. J. Inf. Secur.* **18**, 801–814 (2019).
32. Cartwright, A. J. The elephant in the room: cybersecurity in healthcare. *J. Clin. Monit. Comput.* **37**, 1123–1132 (2023).
33. Prompt Infection: LLM-to-LLM Prompt Injection within Multi-Agent Systems. https://arxiv.org/html/2410.07283v1.
34. Wang, C. et al. Hybrid Alignment Training for Large Language Models. In *Findings of the Association for Computational Linguistics: ACL 2024*, 11389–11403 (Bangkok Thailand, Association for Computational Linguistics, 2024).
35. Research. https://www.anthropic.com/research#alignment.
36. Debenedetti, E. et al. AgentDojo: a dynamic environment to evaluate attacks and defenses for LLM agents. *NeurIPS* (2024).
37. Han, T. et al. Comparative analysis of multimodal large language model performance on clinical vignette questions. *JAMA* **331**, 1320–1321 (2024).
38. Yang, Y., Jin, Q., Huang, F. & Lu, Z. Adversarial attacks on Large Language Models in medicine. *arXiv [cs.AI]* (2024).
39. Lu, A., Zhang, H., Zhang, Y., Wang, X. & Yang, D. Bounding the capabilities of large language models in open text generation with prompt constraints. 1982–2008 (2023).
40. Authorship and AI tools. COPE: Committee on Publication Ethics https://publicationethics.org/cope-position-statements/ai-author.

of the NHS, the NIHR, or the Department of Health and Social Care. This work was funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union. Neither the European Union nor the granting authority can be held responsible for them.

## Competing interests
The authors declare the following competing interests: DT received honoraria for lectures by Bayer and holds shares in StratifAI GmbH, Germany. SF has received honoraria from MSD and BMS. TJB is the owner of Smart Health Heidelberg GmbH, Heidelberg, Germany, outside of the scope of the submitted work. JNK declares consulting services for Bioptimus, France; Owkin, France; DoMore Diagnostics, Norway; Panakeia, UK; AstraZeneca, UK; Mindpeak, Germany; and MultiplexDx, Slovakia. Furthermore, he holds shares in StratifAI GmbH, Germany, Synagen GmbH, Germany, and has received a research grant by GSK, and has received honoraria by AstraZeneca, Bayer, Daiichi Sankyo, Eisai, Janssen, Merck, MSD, BMS, Roche, Pfizer, and Fresenius. ICW has received honoraria from AstraZeneca. DF holds shares in Synagen GmbH, Germany. No other competing interests are declared by any of the remaining authors.

## Additional information
**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s41467-024-55631-x.

**Correspondence** and requests for materials should be addressed to Jakob Nikolas Kather.

**Peer review information** *Nature Communications* thanks Hao Chen, and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

**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.