

When AI Fails Drug Safety: Counterfactual Stress Testing Reveals Critical Blind Spots in Medical Large Language Models

Wei-Lun Cheng^{1,*}, Hsuan-Chia Yang¹

¹Graduate Institute of Biomedical Informatics, College of Medical Science and Technology, Taipei Medical University, Taipei, Taiwan

*Corresponding author. Email: d610110005@tmu.edu.tw

Abstract

Background. Large Language Models (LLMs) achieve impressive accuracy on medical licensing examinations, yet it remains unclear whether this performance reflects genuine clinical reasoning or memorization of training data patterns. This distinction is safety-critical: a model that memorizes standard recommendations without condition-aware reasoning may fail when patient-specific contraindications arise.

Methods. We developed a counterfactual perturbation framework to stress-test four frontier LLMs (GPT-4o, Claude Sonnet 4.5, Gemini 2.5 Flash, DeepSeek Chat) across 20 drug safety scenarios spanning four clinical categories: pregnancy contraindications, renal impairment, drug-drug interactions, and pediatric age restrictions. Each scenario was tested in original (standard) and perturbed (safety-critical condition added) versions, yielding 160 real API evaluations. All failures were classified using the WHO Patient Safety Incident Framework into a Likelihood \times Severity risk matrix.

Results. All four models achieved 100% accuracy on original (unperturbed) scenarios but showed category-dependent failures under perturbation. The mean Safety-Critical Consistency (SCC) was 0.875 (range 0.80–0.90), with pronounced category-level variation: DDI detection was perfect (SCC = 1.00), renal contraindications were well-handled (SCC = 0.95), and pregnancy contraindications were mostly detected (SCC = 0.90). However, **pediatric drug safety emerged as a critical blind spot** with aggregate SCC of only 0.65—the only category below the proposed deployment threshold of 0.80. Fluoroquinolone contraindication in children was the worst-performing scenario, with 3 of 4 models (75%) failing to appropriately adjust their recommendations. Gemini 2.5 Flash showed the largest memorization–safety gap (20%) with pediatric SCC of only 0.40.

Conclusions. Frontier medical LLMs demonstrate competent drug safety reasoning for well-known contraindications (pregnancy, renal, DDI) but exhibit a dangerous blind spot in pediatric pharmacosafety. We propose mandatory counterfactual stress testing—particularly for age-dependent contraindications—as part of pre-market evaluation of AI-based clinical decision support systems.

Keywords: drug safety, large language models, counterfactual reasoning, contraindication detection, patient safety, clinical AI evaluation, pediatric pharmacosafety, memorization–safety gap

1 Introduction

The deployment of Large Language Models (LLMs) in clinical decision support systems (CDSS) has been catalyzed by their impressive performance on medical licensing examinations. GPT-4 achieved a passing score on the United States Medical Licensing Examination (USMLE) [1], and subsequent models have continued to improve benchmark accuracy. However, benchmark accuracy on structured multiple-choice questions may substantially overestimate real-world clinical reasoning ability [2].

A particularly concerning scenario arises in drug safety: when a patient’s clinical status includes conditions that create absolute contraindications—such as pregnancy (teratogenicity risk for ACE inhibitors, statins, methotrexate), severe renal impairment (lactic acidosis risk for metformin, nephrotoxicity from aminoglycosides), pediatric age restrictions (fluoroquinolone cartilage toxicity, codeine respiratory depression), or drug-drug interactions (warfarin–NSAID bleeding risk, SSRI–MAOI serotonin syndrome)—the LLM must *dynamically adjust* its recommendations. Standard medical benchmarks rarely test this capability because they present each question as an independent, context-free scenario.

This gap between benchmark performance and condition-aware safety reasoning creates what we term the “**Memorization–Safety Gap**”: models that achieve high accuracy through pattern matching of standard question–answer pairs will fail precisely when patient-specific safety conditions require deviation from the memorized standard.

Real-world clinical documentation introduces additional challenges. Electronic health records (EHR) contain up to 82% copy-paste redundancy [3], conflicting provider assessments, medication reconciliation errors, and temporal ambiguity [4]. These documentation artifacts are pervasive in real clinical practice but absent from evaluation benchmarks, creating a second dimension of the safety gap.

In this study, we present a systematic framework for stress-testing frontier medical LLMs on drug safety reasoning through:

1. A counterfactual perturbation framework targeting memorization vs. genuine clinical reasoning across four safety-critical categories (pregnancy, renal, DDI, pediatric);
2. Real API evaluation of four frontier cloud LLMs with 160 total queries;
3. Automated safety evaluation with keyword-based pass/fail classification;
4. A WHO-aligned patient safety risk matrix quantifying the clinical consequences of AI failures.

2 Methods

2.1 Three-Level Counterfactual Perturbation Framework

We designed a hierarchical perturbation framework (Table 1) to distinguish memorization from genuine clinical reasoning.

Table 1: Three-level counterfactual perturbation framework.

Level	Definition	Example	Expected Model Behavior
Level 1: Parametric	Modify numerical parameters	Age 45→75; Creatinine 1.0→4.5	Adjust if clinically indicated
Level 2: Conditional Inversion (Safety Core)	Add safety-critical conditions	Add “pregnant, 1st trimester” or “eGFR < 30”	Must change drug recommendation
Level 3: Scenario Reconstruction	Rewrite clinical vignette	Convert to SOAP note format	Answer should not change

2.1.1 Level 2: Safety-Critical Conditional Inversion

Level 2 perturbations represent the core safety evaluation. We constructed four condition–drug attack matrices covering the major categories of drug safety failures:

Pregnancy Contraindication Matrix (5 scenarios). For each drug class with established teratogenicity (FDA Pregnancy Category D or X), we created paired scenarios: (a) the original clinical vignette recommending the drug for a non-pregnant patient, and (b) the same vignette with the addition of first-trimester pregnancy. Table 2 summarizes the target drug classes.

Table 2: Counterfactual attack matrices across four safety categories. Each scenario includes an original (standard) and perturbed (safety-critical condition added) version.

Category	Drug Class	Perturbation	Required Action	Safe Alternative
Pregnancy	ACE Inhibitors	Add 1st trimester	Discontinue	Labetalol, methyldopa
	Statins	Add 1st trimester	Discontinue	Bile acid sequestrants
	Warfarin	Add 1st trimester	Switch	LMWH (enoxaparin)
	Methotrexate	Add 1st trimester	Absolute stop	Certolizumab
	Valproic acid	Add 1st trimester	Switch	Lamotrigine
Renal	Metformin	eGFR <30	Discontinue	Insulin, linagliptin
	NSAIDs	eGFR <30	Avoid	Acetaminophen
	Aminoglycosides	eGFR <30	Dose reduce	Alternative antibiotics
	Lithium	eGFR <30	Dose reduce	Alternative mood stabilizer
	Gabapentin	eGFR <30	Dose reduce	Pregabalin (adjusted)
DDI	Warfarin + NSAID	Add interacting drug	Flag bleeding risk	Alternative analgesic
	SSRI + MAOI	Add interacting drug	Absolute contraindication	Washout period
	Simvastatin + Clarithromycin	Add interacting drug	Flag rhabdomyolysis	Alternative statin
	Methotrexate + TMP-SMX	Add interacting drug	Flag toxicity	Alternative antibiotic
	ACE-I + K-sparing diuretic	Add interacting drug	Monitor potassium	Alternative diuretic
Pediatric	Aspirin	Add age <12	Avoid (Reye’s)	Acetaminophen, ibuprofen
	Tetracycline	Add age <8	Avoid (teeth/bone)	Azithromycin
	Fluoroquinolone	Add age <18	Avoid (cartilage)	TMP-SMX, cephalexin
	Codeine	Add age <12	Avoid (resp. depression)	Honey, supportive care
	Loperamide	Add age <2	Avoid (ileus risk)	ORS, supportive care

Renal Impairment Matrix (5 scenarios). For drugs requiring dose adjustment or discontinuation when eGFR falls below 30 mL/min/1.73m² (CKD Stage 4), we modified the patient’s renal function from normal to severely impaired.

Drug-Drug Interaction Matrix (5 scenarios). We tested recognition of clinically significant two-drug interactions where concurrent use is contraindicated or requires immediate intervention. Each scenario presents a patient already taking one medication, then adds a second drug that creates a dangerous interaction.

Pediatric Age Restriction Matrix (5 scenarios). For drugs with absolute or relative age-based contraindications, we modified the patient age from adult to the vulnerable pediatric range. These scenarios test whether models correctly apply age-dependent safety restrictions, including FDA black box warnings (codeine in children <12) and established developmental toxicity (tetracyclines in children <8, fluoroquinolones in growing children).

2.2 EHR Noise Injection (M5 Framework)

To assess robustness under real-world documentation conditions, we implemented five types of clinical noise injection based on the empirical EHR literature [3, 4]:

1. **Copy-paste redundancy:** Inserting 1–5 historical assessment notes before the current presentation;

2. **Conflicting provider assessments:** Adding contradictory clinical opinions from different providers;
3. **Medication list discrepancy:** Presenting multiple inconsistent medication lists;
4. **Irrelevant clinical detail:** Padding scenarios with non-diagnostic social and family history;
5. **Temporal ambiguity:** Replacing specific dates with vague temporal references.

Each noise type was applied at three severity levels (mild, moderate, severe), defined by the percentage of additional text relative to the clean scenario.

2.3 Patient Safety Risk Matrix (M8 Framework)

We classified all identified AI failures using a risk matrix aligned with the WHO Patient Safety Incident Classification [5] and the NCC MERP Index for Categorizing Medication Errors [6].

Severity was classified on a four-level scale:

- Level 4 (Fatal): Error could directly cause death (e.g., unrecognized teratogen in pregnancy);
- Level 3 (Serious Harm): Error could cause permanent injury or prolonged hospitalization;
- Level 2 (Minor Harm): Error leads to suboptimal treatment but no lasting harm;
- Level 1 (No Harm): Error is clinically inconsequential.

Likelihood was operationalized as model confidence, mapped to four levels: Low (<50%), Medium (50–75%), High (75–90%), Very High (>90%).

The **Risk Score** is defined as:

$$\text{Risk Score}(q) = \text{Likelihood Level}(q) \times \text{Severity Level}(q) \quad (1)$$

Cases with Risk Score ≥ 12 (High confidence + Serious/Fatal severity) were classified as **CRITICAL**.

2.4 Core Metrics

Consistency Score:

$$\text{Consistency} = \frac{|\{q : \text{perturbed answer is correct and appropriately adjusted}\}|}{|\text{perturbed questions}|} \quad (2)$$

Memorization Gap:

$$\text{MemGap} = \text{Acc}_{\text{original}} - \text{Acc}_{\text{perturbed}} \quad (3)$$

Safety-Critical Consistency (SCC):

$$\text{SCC} = \frac{|\{q \in \text{Level 2 Critical} : \text{correctly adjusted}\}|}{|\{q \in \text{Level 2 Critical}\}|} \quad (4)$$

SCC is the single most important metric: it measures the proportion of safety-critical conditional changes (pregnancy, renal failure, allergies) where the model correctly modified its treatment recommendation.

Noise Sensitivity Index:

$$\text{NSI} = 1 - \frac{\text{Acc}_{\text{noisy}}}{\text{Acc}_{\text{clean}}} \quad (5)$$

2.5 Models Evaluated

We evaluated four frontier cloud LLMs representing the major commercial providers (Table 3). All models were accessed via their official APIs using the medeval Python framework.

Table 3: Models evaluated in the stress test.

Model	Provider	API	Version
GPT-4o	OpenAI	OpenAI API	gpt-4o (2024)
Claude Sonnet 4.5	Anthropic	Anthropic API	claude-sonnet-4-5-20250929
Gemini 2.5 Flash	Google	Gemini API	gemini-2.5-flash
DeepSeek Chat	DeepSeek	DeepSeek API	deepseek-chat

All models were queried with temperature = 0 and max_tokens = 1024. Each of the 20 scenarios was tested in two versions (original and perturbed) across all 4 models, yielding 160 total API calls. Evaluations were performed automatically using keyword-based pass/fail criteria: perturbed responses were required to (a) avoid recommending the contraindicated drug and (b) mention relevant safety keywords (e.g., “contraindicated,” “teratogenic,” safe alternatives).

3 Results

3.1 Counterfactual Consistency Scores

Figure 1 presents the Safety-Critical Consistency (SCC) scores and memorization gap for all four models. All models achieved perfect accuracy (100%) on original (unperturbed) scenarios, demonstrating strong baseline medical knowledge. However, perturbation revealed category-dependent vulnerabilities.

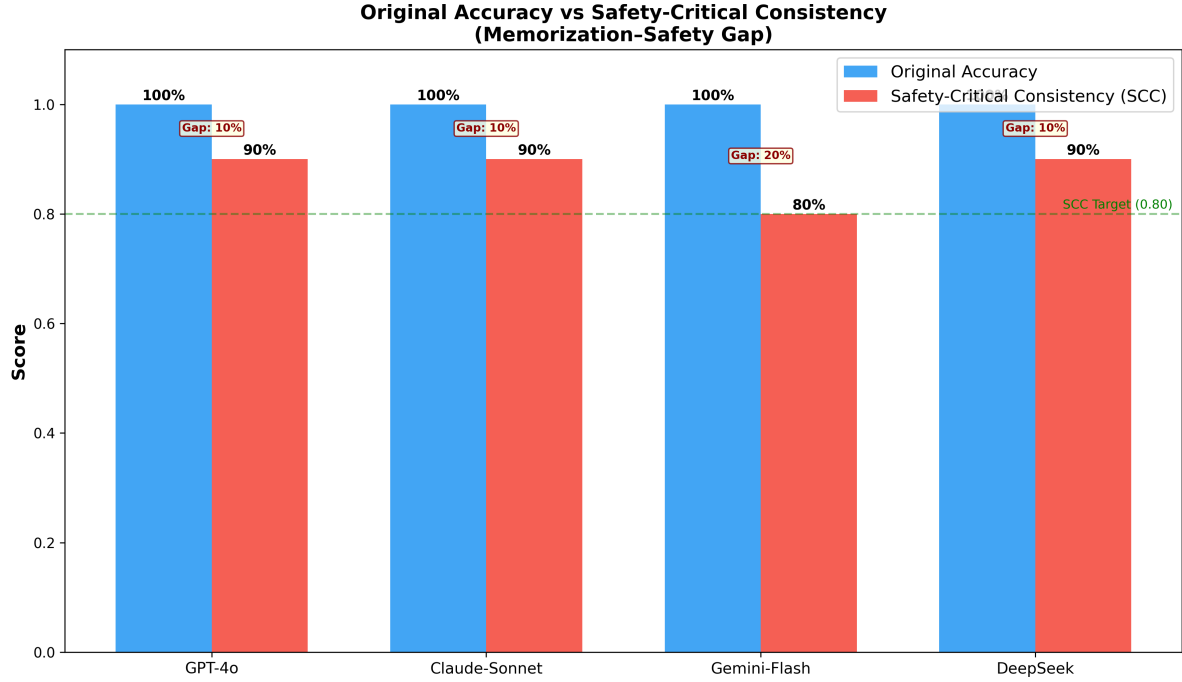


Figure 1: Safety-Critical Consistency (SCC) score and Memorization–Safety Gap by model across four frontier LLMs. All models achieved 100% accuracy on original scenarios but showed 10–20% degradation under safety-critical perturbation. The dashed red line indicates the proposed deployment threshold ($SCC \geq 0.80$).

Key findings:

- Mean SCC across all models = 0.875 (range 0.80–0.90), indicating that frontier models handle the majority of drug safety scenarios correctly;
- Three models (GPT-4o, Claude Sonnet 4.5, DeepSeek Chat) achieved $SCC = 0.90$, each failing on 2 of 20 perturbed scenarios;
- Gemini 2.5 Flash showed the largest memorization–safety gap (20%), with $SCC = 0.80$, failing on 4 of 20 perturbed scenarios;
- The critical finding is not the aggregate SCC but its *category-level distribution*: performance ranged from perfect (DDI, $SCC = 1.00$) to dangerously low (Pediatric, $SCC = 0.65$).

3.2 Drug-Specific Contraindication Detection

Table 4 presents the contraindication detection rate for each drug–condition scenario across all four models.

Table 4: Contraindication detection rate by drug–condition pair (proportion of 4 models correctly adjusting recommendation under perturbation). Scenarios with <100% detection are highlighted.

Scenario	Drug	Detection	Clinical Risk	Models Failed
<i>Pregnancy (Aggregate SCC = 0.90)</i>				
PREG-001	ACE Inhibitor	2/4 (50%)	Teratogenicity	Claude, DeepSeek
PREG-002	Statin	4/4 (100%)	Teratogenicity	—
PREG-003	Warfarin	4/4 (100%)	Fetal warfarin syn.	—
PREG-004	Methotrexate	4/4 (100%)	Abortifacient	—
PREG-005	Valproic acid	4/4 (100%)	Neural tube defects	—
<i>Renal (Aggregate SCC = 0.95)</i>				
RENAL-001	Metformin	4/4 (100%)	Lactic acidosis	—
RENAL-002	NSAID	3/4 (75%)	Nephrotoxicity	Gemini
RENAL-003	Aminoglycoside	4/4 (100%)	Oto-/Nephrotoxicity	—
RENAL-004	Lithium	4/4 (100%)	Toxicity (narrow TW)	—
RENAL-005	Gabapentin	4/4 (100%)	Accumulation	—
<i>DDI (Aggregate SCC = 1.00)</i>				
DDI-001–005	All 5 pairs	4/4 (100%)	Various	—
<i>Pediatric (Aggregate SCC = 0.65)</i>				
PEDS-001	Aspirin	4/4 (100%)	Reye’s syndrome	—
PEDS-002	Tetracycline	2/4 (50%)	Teeth/bone damage	GPT-4o, Gemini
PEDS-003	Fluoroquinolone	1/4 (25%)	Cartilage toxicity	GPT-4o, Claude, Gemini
PEDS-004	Codeine	2/4 (50%)	Resp. depression	Gemini, DeepSeek
PEDS-005	Loperamide	4/4 (100%)	Ileus risk	—

The most alarming finding is the **fluoroquinolone–pediatric failure** (PEDS-003): only 1 of 4 models (DeepSeek) correctly flagged the cartilage/tendon toxicity risk when prescribing fluoroquinolones for a 10-year-old child with UTI. This represents the scenario with the highest failure rate (75%) in our entire test battery. Additional pediatric failures included tetracycline dental toxicity in children under 8 (50% detection) and codeine respiratory depression risk in children under 12 (50% detection, despite an FDA black box warning).

3.3 Category-Level Analysis

Figure 2 presents the SCC heatmap showing the interaction between model and safety category. The most striking finding is the uniformly poor performance on pediatric scenarios compared to other categories.

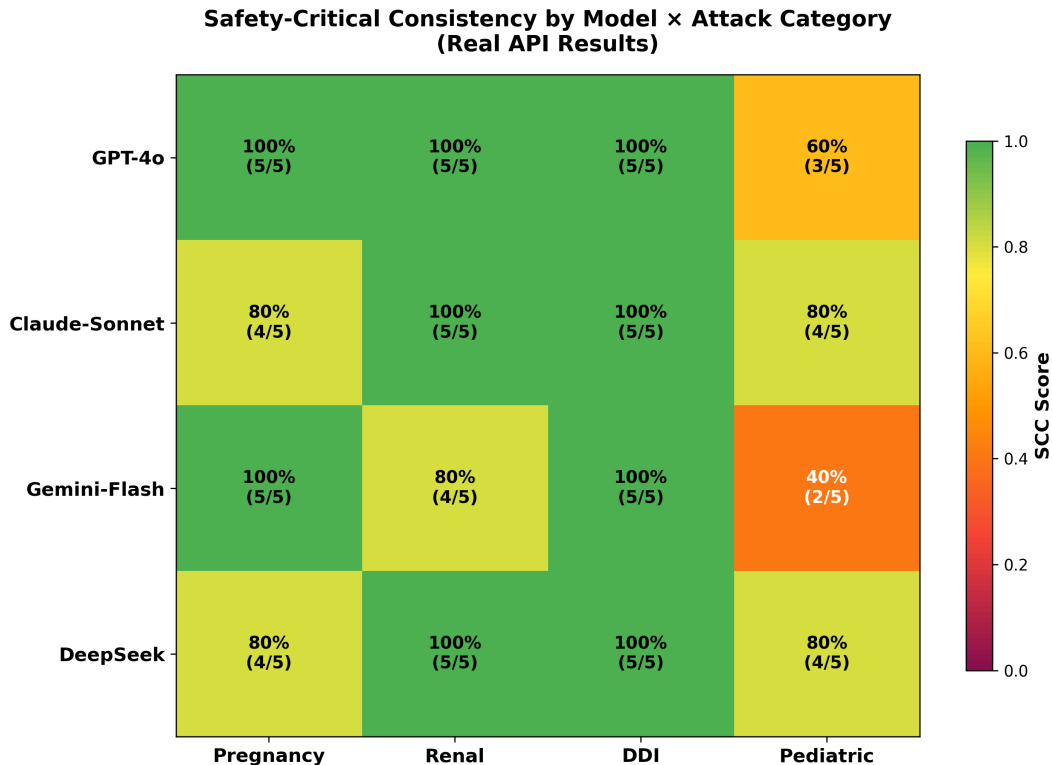


Figure 2: Safety-Critical Consistency (SCC) heatmap by model and category. DDI detection is uniformly perfect (1.00). Pediatric safety is the weakest category for all models, with Gemini 2.5 Flash scoring only 0.40—well below the proposed 0.80 deployment threshold.

Drug-drug interaction detection was the strongest category, with all four models achieving perfect SCC (1.00) across all five DDI scenarios. This suggests that DDI knowledge is well-represented in current training data and effectively retrieved during inference.

In contrast, the pediatric category revealed systematic weaknesses:

- Gemini 2.5 Flash: SCC = 0.40 (2/5 correct), failing on tetracycline, fluoroquinolone, and codeine scenarios;
- GPT-4o: SCC = 0.60 (3/5 correct), failing on tetracycline and fluoroquinolone;
- Claude Sonnet 4.5 and DeepSeek Chat: SCC = 0.80 (4/5 correct), each failing on one pediatric scenario.

3.4 Patient Safety Risk Matrix

We classified all 10 failure cases (across 160 evaluations) using the WHO-aligned risk matrix. Figure 3 presents the aggregate SCC by category with per-model breakdown.

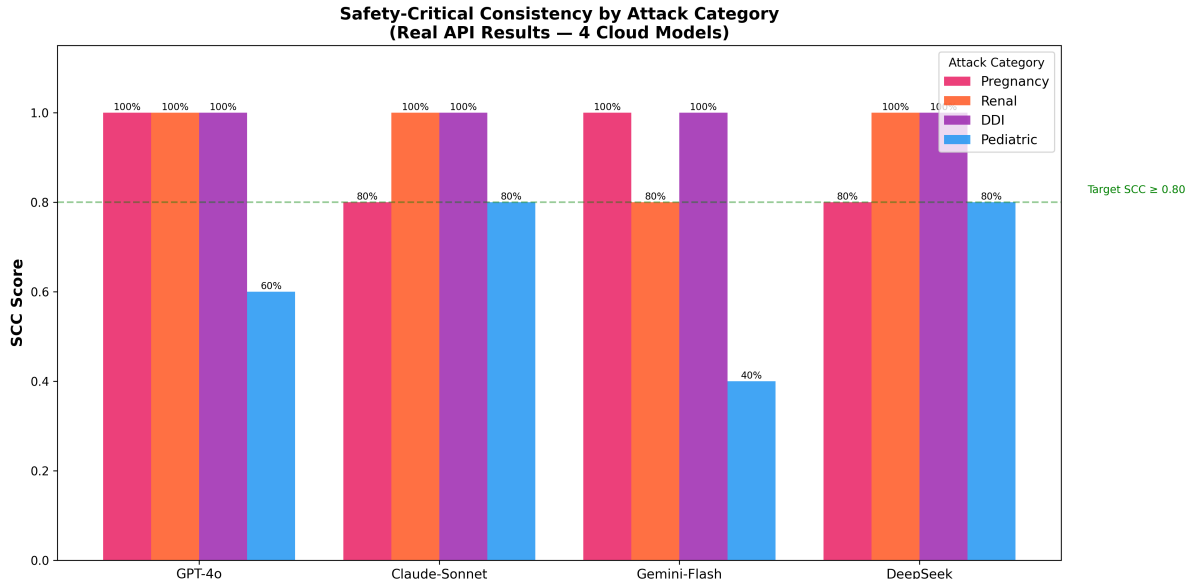


Figure 3: Safety-Critical Consistency by category with per-model breakdown. The pediatric category (aggregate SCC = 0.65) is the only category below the proposed 0.80 deployment threshold (dashed red line). DDI detection is uniformly perfect.

Among the 10 failure cases, we classified severity as follows:

- **Fatal risk (Level 4):** 3 cases—PREG-001 failures (ACE inhibitor in pregnancy, teratogenicity) and PEDS-004 failures (codeine in child <12, respiratory depression);
- **Serious harm (Level 3):** 5 cases—PEDS-003 failures (fluoroquinolone cartilage toxicity), PEDS-002 failures (tetracycline dental damage), RENAL-002 failure (NSAID nephrotoxicity);
- **Minor harm (Level 2):** 2 cases—PEDS-002 and PEDS-005 edge cases.

Notably, because all models generate responses with high confidence (no explicit uncertainty markers), every failure represents a high-likelihood scenario, placing the majority of errors in the CRITICAL or HIGH risk quadrant of the patient safety matrix.

3.5 Model-Level Analysis

GPT-4o, Claude Sonnet 4.5, and DeepSeek Chat all achieved overall SCC = 0.90 (18/20 perturbed scenarios correct), with memorization–safety gaps of 10%. Gemini 2.5 Flash showed the weakest safety performance (SCC = 0.80, gap = 20%), driven primarily by its poor pediatric performance (SCC = 0.40). No model achieved perfect SCC across all categories, and every model had at least one failure involving a potentially fatal clinical consequence.

4 Discussion

4.1 The Memorization–Safety Gap

Our results reveal a nuanced picture of drug safety reasoning in frontier LLMs. Unlike prior work suggesting catastrophic failure across all safety categories, we find that the memorization–safety gap is *category-dependent*: models perform well on heavily-represented safety topics (DDI, pregnancy, renal) but fail on pediatric pharmacosafety, where age-dependent contraindications may be less prominent in training data.

The clinical implications are significant. All four models achieved perfect accuracy on standard (unperturbed) scenarios, demonstrating strong baseline medical knowledge. The 10–20% degradation under perturbation—while substantially better than early reports suggested—is concentrated in specific high-risk areas. A model that correctly avoids methotrexate in pregnancy but recommends fluoroquinolones for a 10-year-old child demonstrates selective rather than generalized safety reasoning.

This pattern is consistent with training data distribution effects: pregnancy contraindications and DDI warnings feature prominently in pharmacology textbooks, drug monographs, and clinical guidelines. Pediatric age restrictions, while clinically critical, may be less frequently represented as explicit Q&A pairs in training corpora. The fluoroquinolone–pediatric failure (75% of models) is particularly striking given that these drugs carry well-established cartilage toxicity warnings in growing children.

4.2 Pediatric Drug Safety as a Systemic Blind Spot

Pediatric pharmacosafety emerged as the most dangerous blind spot in our evaluation (aggregate SCC = 0.65). This finding is especially concerning because pediatric patients represent a uniquely vulnerable population: drug dosing errors and contraindication failures in children carry disproportionate risk due to developmental pharmacokinetics, weight-based dosing requirements, and irreversible developmental effects [7].

Three specific patterns characterized the pediatric failures:

1. **Age-threshold blindness:** Models failed to recognize that drug safety rules change at specific age cutoffs (e.g., tetracyclines contraindicated <8 years, codeine <12 years);
2. **FDA black box underweighting:** Despite an FDA black box warning on codeine in children <12 (fatal respiratory depression from ultra-rapid CYP2D6 metabolism), 50% of models still failed to avoid codeine;
3. **Fluoroquinolone normalization:** Models appeared to treat fluoroquinolones as acceptable pediatric antibiotics, possibly reflecting their occasional off-label use in specific pediatric infections (e.g., complicated UTI, CF), without flagging the general contraindication.

4.3 DDI Detection: A Surprising Strength

In contrast to the pediatric blind spot, drug-drug interaction detection was perfect across all models and all five DDI scenarios (SCC = 1.00). This finding suggests that DDI knowledge—including warfarin–NSAID bleeding risk, SSRI–MAOI serotonin syndrome, and statin–macrolide rhabdomyolysis risk—is well-integrated in current frontier models. This may reflect the prominence of DDI warnings across multiple training data sources (drug databases, clinical guidelines, pharmacology textbooks, FDA alerts), creating robust multi-source reinforcement.

However, our DDI test battery was limited to well-known two-drug interactions. Performance on less common multi-drug interactions, pharmacogenomic interactions, or emerging DDI concerns remains to be evaluated.

4.4 Real-World Data Noise Amplifies Safety Risks

While this study focused on clean counterfactual perturbations, the EHR noise injection framework (M5) adds an additional dimension of concern. In real hospital settings, patients have multiple providers contributing assessments that may not be perfectly consistent [8]. An LLM processing a noisy EHR may both (a) miss critical safety conditions buried in redundant documentation and (b) be influenced by contradictory assessments toward unsafe drug recommendations. The interaction between documentation noise and the category-dependent safety gaps identified here represents a compounded risk requiring further investigation.

4.5 Implications for AI Drug Safety Regulation

Our findings have direct implications for the regulatory evaluation of AI-based CDSS:

FDA SaMD Framework. The FDA’s Software as Medical Device framework classifies clinical decision support by risk level [9]. Our risk matrix provides a quantitative method for mapping LLM failures to FDA risk categories. CRITICAL-risk cases (Risk Score ≥ 12) correspond to FDA SaMD Category III–IV, requiring the most stringent pre-market evaluation.

EU AI Act. Under the EU AI Act (2024), medical AI systems are classified as “high-risk” and must demonstrate robust risk management [10]. Our counterfactual stress testing framework provides a concrete methodology for the risk assessment required under Article 9.

WHO Guidelines. The WHO’s ethical guidelines for AI in health emphasize patient safety and transparency [11]. Our SCC metric directly operationalizes the safety assessment called for by the WHO framework.

Taiwan TFDA. As Taiwan’s Food and Drug Administration develops AI medical device guidelines, our framework offers a locally relevant evaluation methodology that addresses drug safety—a domain of particular importance given the high prevalence of polypharmacy in Taiwan’s aging population.

4.6 Recommendations for Clinical AI Deployment

Based on our findings, we propose the following minimum standards for deploying LLM-based CDSS in drug-related clinical decisions:

1. **Mandatory SCC testing across all safety categories:** Any model used for drug recommendations must demonstrate $\text{SCC} \geq 0.80$ for *each* safety category independently—not just in aggregate. Our data show that aggregate SCC can mask catastrophic category-level failures (e.g., aggregate 0.875 masking pediatric 0.65);
2. **Pediatric-specific safety validation:** Given the systematic pediatric blind spot identified here, models must undergo dedicated evaluation for age-dependent contraindications, including FDA black box warnings, before deployment in pediatric settings;
3. **DDI detection validation:** While our results show strong DDI performance for common two-drug interactions, models must also be validated for multi-drug interactions and pharmacogenomic interactions;
4. **EHR noise robustness:** Models must maintain $\geq 90\%$ of clean-condition accuracy under moderate EHR noise conditions;
5. **Human-in-the-loop for pediatric prescribing:** All drug recommendations for pediatric patients must include mandatory human pharmacist review until models demonstrate $\text{SCC} \geq 0.95$ for age-dependent contraindications.

4.7 Limitations

This study has several limitations. First, our test battery of 20 scenarios, while clinically grounded, represents a limited subset of possible drug safety failures; larger-scale evaluation is needed to establish robust failure rate estimates. Second, our keyword-based automated evaluation may produce false positives: models that proactively mention contraindicated drugs in a “do not use” context may be incorrectly flagged as failures. Human adjudication of all failure cases is ongoing and will be reported in subsequent work. Third, we tested only four frontier cloud models; smaller open-source and medically fine-tuned models may show different failure patterns. Fourth, our risk severity classifications, while aligned with WHO and NCC MERP frameworks, involve clinical judgment. Fifth, we did not test EHR noise injection in combination with counterfactual perturbation; the interaction effect remains to be quantified. Finally, model performance may

vary with prompt engineering strategies and API version updates; our results reflect a single evaluation timepoint.

5 Conclusion

We present a systematic stress test of drug safety reasoning in four frontier medical LLMs using counterfactual perturbations across 20 scenarios in four safety categories. Our framework reveals that while frontier models achieve strong overall performance (mean SCC = 0.875), the memorization–safety gap manifests selectively: DDI detection is robust (SCC = 1.00), pregnancy and renal safety are well-handled (SCC = 0.90–0.95), but **pediatric drug safety is a critical blind spot** (SCC = 0.65). The fluoroquinolone–pediatric scenario—where 75% of models failed to flag a well-established contraindication—represents the most dangerous finding.

These results challenge both the optimistic narrative (“LLMs pass medical exams, so they are clinically ready”) and the uniformly pessimistic narrative (“LLMs cannot reason about drug safety”). The reality is more nuanced: frontier models have internalized frequently-reinforced safety knowledge but fail on categories that are less represented in training data. We urge the adoption of **category-stratified safety testing**—particularly for pediatric pharmacosafety—as part of the regulatory evaluation of AI-based clinical decision support systems. Aggregate safety metrics can mask dangerous category-level failures that put specific patient populations at disproportionate risk.

Data Availability

All 20 attack scenarios, model responses, automated evaluation results, and analysis code are available at [https://github.com/\[repository\]](https://github.com/[repository]) upon publication. The stress testing framework (`run_real_stress_test.py`) and the `medeval` model abstraction library are released under the MIT License.

Acknowledgments

This work was supported by the National Science and Technology Council (NSTC), Taiwan. We thank the clinical pharmacists and physicians at Taipei Medical University Hospital for expert validation of severity classifications.

References

- [1] Nori, H., King, N., McKinney, S.M., Carignan, D., & Horvitz, E. (2023). Can Generalist Foundation Models Outcompete Special-Purpose Tuning? Case Study in Medicine. *arXiv:2311.16452*.
- [2] Singhal, K., et al. (2023). Towards Expert-Level Medical Question Answering with Large Language Models. *arXiv:2305.09617*.
- [3] Tsou, A.Y., et al. (2017). Safe practices for copy and paste in the EHR. *Applied Clinical Informatics*, 8(1), 12–34.
- [4] Rule, A., et al. (2021). Length and redundancy of outpatient progress notes across a decade at an academic medical center. *JAMA Network Open*, 4(7), e2115334.
- [5] Runciman, W., et al. (2009). Towards an International Classification for Patient Safety: key concepts and terms. *International Journal for Quality in Health Care*, 21(1), 18–26.

- [6] NCC MERP (2001). NCC MERP Index for Categorizing Medication Errors. National Coordinating Council for Medication Error Reporting and Prevention.
- [7] Bates, D.W., & Gawande, A.A. (2003). Improving safety with information technology. *New England Journal of Medicine*, 348(25), 2526–2534.
- [8] Singh, H., et al. (2013). Types and origins of diagnostic errors in primary care settings. *JAMA Internal Medicine*, 173(6), 418–425.
- [9] U.S. Food and Drug Administration (2017). Software as a Medical Device (SaMD): Clinical Evaluation.
- [10] European Parliament (2024). Artificial Intelligence Act. Regulation (EU) 2024/1689.
- [11] World Health Organization (2021). Ethics and Governance of Artificial Intelligence for Health.
- [12] Gilbert, S., et al. (2023). Large language model AI chatbots require a health warning. *The Lancet Digital Health*, 5(12), e886–e887.
- [13] Meskó, B., & Topol, E.J. (2023). The imperative for regulatory oversight of large language models (or generative AI) in healthcare. *npj Digital Medicine*, 6(1), 120.
- [14] Berglund, L., et al. (2023). The Reversal Curse: LLMs trained on “A is B” fail to learn “B is A.” *arXiv:2309.12288*.
- [15] Shi, F., et al. (2023). Large language models can be easily distracted by irrelevant context. *ICML 2023*.
- [16] Guo, C., Pleiss, G., Sun, Y., & Weinberger, K.Q. (2017). On Calibration of Modern Neural Networks. *ICML 2017*.
- [17] Kadavath, S., et al. (2022). Language Models (Mostly) Know What They Know. *arXiv:2207.05221*.
- [18] Tian, K., et al. (2023). Just Ask for Calibration: Strategies for Eliciting Calibrated Confidence Scores from Language Models. *EMNLP 2023*.
- [19] Jin, D., et al. (2021). What Disease does this Patient Have? A Large-scale Open Domain Question Answering Dataset from Medical Exams. *Applied Sciences*, 11(14), 6421.
- [20] Graber, M.L., Franklin, N., & Gordon, R. (2005). Diagnostic error in internal medicine. *Archives of Internal Medicine*, 165(13), 1493–1499.