

Payam Hosseinzadeh Kasani^{1,2}, Kee Hyun Cho^{1,2}, Cheol-Heui Yun^{3,4,5*}

¹ Department of Pediatrics, Kangwon National University Hospital, Chuncheon, Republic of Korea. ² Department of Pediatrics, Kangwon National University School of Medicine, Chuncheon, Republic of Korea. ³ Department of Agricultural Biotechnology, and Research Institute of Agriculture and Life Sciences, Seoul National University, Seoul, Republic of Korea. ⁴ Center for Food and Bioconvergence, and Interdisciplinary Programs in Agricultural Genomics, Seoul National University, Seoul, Republic of Korea. ⁵ Institutes of Green Bio Science and Technology, Seoul National University, Pyeongchang, Republic of Korea.

Abstract

Large language models are rapidly entering medicine for documentation, patient communication, education, administration, research synthesis, and decision support, but their fluent outputs can be mistaken for trustworthy medical knowledge.

This review examines the limitations of large language models in medicine across three domains: first, epistemic reliability, which includes hallucinations, inconsistent outputs, outdated knowledge, and poor uncertainty calibration; second, clinical translation, which includes workflow disruption, diagnostic risk, unsafe patient-facing use, and overreliance; and third, governance accountability, which includes transparency, bias, privacy, cybersecurity, liability, and regulation.

Although LLMs may support selected clinical tasks and diagnostic second opinions, their context-dependent performance and hidden failure modes require local validation, human oversight, privacy protection, subgroup testing, auditability, and continuous governance to ensure they remain assistive rather than autonomous clinical tools.

Determine whether fluent output reflects trustworthy medical knowledge

- Fluency- trustworthiness gap
- Hallucination and misinformation
- Risk of Overinterpreting polished answers



Assess consistency, reproducibility, and knowledge currency

- Factual instability across prompts and settings
- Non-reproducibility across model versions
- Knowledge obsolescence and guideline drift

Evaluate real-world clinical reasoning rather than benchmark performance alone

- Benchmark-practice mismatch
- Inflexible reasoning and overconfidence
- Need for realistic open-ended case evaluation

Protect ethical, contextual, and patient-centered judgment

- Ethical and contextual reasoning blind spots
- Bias-sensitive and value-laden decision
- High-stakes outputs require clinician oversight

Figure 1. Epistemic reliability limitations of large language models in medicine

01. Bias and Unequal Subgroup Performance

Training corpora may encode inequities in medical literature, clinical documentation, internet data, and health systems; average accuracy may hide subgroup errors.

02. Global Health and Language Inequity

Dominance of English-language and high-resource evidence may produce recommendations that do not match local disease burden, formularies, infrastructure, or health literacy.

03. Privacy, Confidentiality, and Cybersecurity

Prompts, outputs, metadata, or logs may contain PHI and may be processed, retained, reused, or exposed by poorly governed systems; integrated LLM systems may also face prompt injection, data poisoning, or information leakage.

05. Transparency and Explainability Limits

Proprietary models may obscure training data, updates, validation boundaries, and subgroup performance; generated explanations may be post hoc rather than causal.

04. Accountability and Liability Diffusion

Clinical influence may be distributed across clinicians, hospitals, vendors, developers, data providers, and regulators without clear responsibility after harm.

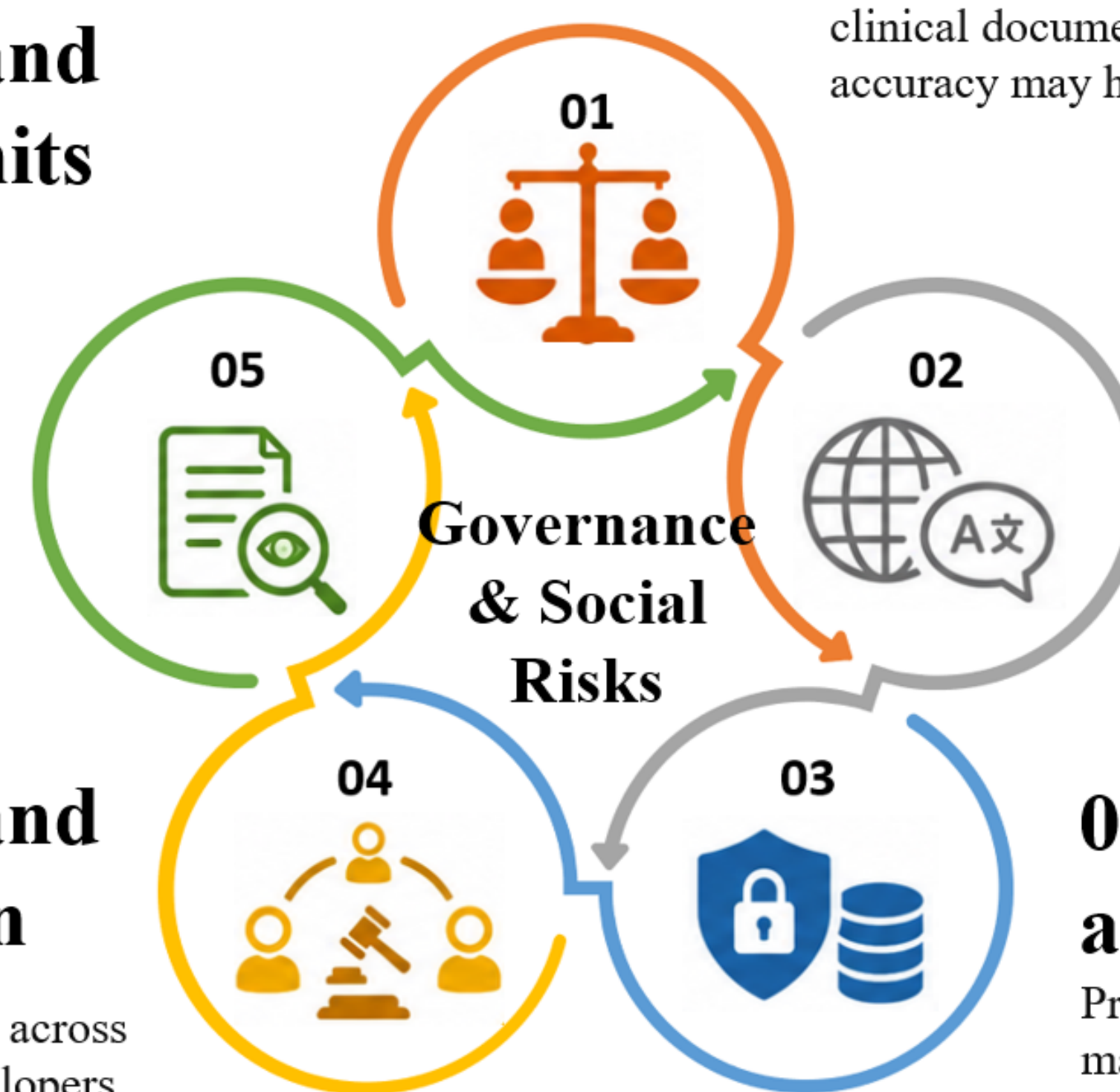


Figure 2. Trust at the System Level: Governance Risks in Medical LLM Deployment

Workflow misalignment and task creep

LLMs may be added to EHR, messaging, documentation, or triage workflows without clear boundaries, routing rules, or escalation criteria, leading to disrupted workflow, duplicated work, and unsafe routing of clinical questions.

Automation Bias and Human Oversight

Fluent, structured outputs may appear authoritative and reduce independent checking, especially among time-pressured or less experienced users, creating overreliance and inappropriate substitution for human judgment.

Quality Assurance and Staged Validation

Strong benchmark or second-opinion performance may not reflect real-world behavior, user reliance, workflow effects, or outcome impact; staged evaluation is required before broad deployment.



Verification Burden and Documentation Propagation

Although drafting is automated, clinicians must still verify accuracy, completeness, tone, medicolegal adequacy, and consistency with the source record; subtle omissions may propagate into later care.

Diagnostic Safety and Patient-Facing Risk

LLMs may generate broad but poorly prioritized differentials or patient messages without adequate red-flag recognition, risk stratification, or uncertainty handling, increasing the risk of missed emergencies and unsafe triage.

Generalizability, Local Adaptation, and Performance Decay

Performance varies by specialty, language, EHR, formulary, local guideline, resource setting, patient subgroup, model update, and prompt template, requiring local validation and repeated re-evaluation.

Figure 3. From Model Performance to Safe Practice: Major Clinical Translation Challenges for Medical LLMs

Table 1. Risk-based governance architecture for medical LLMs

Governance level	Core oversight functions	Required documentation or evidence
Use-case risk classification	Classify each application by intended use, autonomy, patient-facing status, clinical consequence, and user expertise.	Risk tier, intended-use statement, workflow map, prohibited uses, and approval decision.
Developer and vendor oversight	Assess safety, bias, privacy, cybersecurity, data handling, update policy, performance claims, and incident reporting obligations.	Model documentation, validation summary, security report, data-processing agreement, update log, and known limitations.
Institutional implementation	Approve local use cases, integrate with EHR safely, define prompts/templates, train users, and control access.	Local validation report, user-training records, prompt/version registry, privacy-impact assessment, and access-control plan.
Clinical human oversight	Verify outputs, escalate uncertainty, avoid unsupported high-risk use, maintain human judgment, and disclose AI involvement when appropriate.	Clinician review record, escalation pathway, patient-facing disclosure policy, and incident documentation.
Continuous monitoring and AI-QI	Monitor performance decay, subgroup errors, harmful omissions, drift, user behavior, and adverse events; revalidate after updates.	Audit dashboards, periodic performance reports, revalidation schedule, incident-review log, and stopping/modification criteria.
Regulatory, ethics, and publication oversight	Clarify medical-device or clinical decision-support status, define validation standards, require reporting, and support post-market surveillance.	Regulatory classification, ethics checklist, protocol approval, model/prompt disclosure, safeguards, and post-market monitoring plan.

Conclusion

The central challenge of medical LLMs is not their ability to produce fluent text, but whether that text can be trusted, safely translated into clinical workflows, and governed with clear accountability. This review therefore frames LLM deployment as a sociotechnical problem spanning epistemic reliability, clinical implementation, patient safety, equity, privacy, cybersecurity, transparency, and regulation. Ultimately, LLMs should be adopted only as assistive tools under local validation, human oversight, continuous monitoring, and institution-level governance, so that innovation improves care without weakening clinical responsibility or public trust