

Clinical Reasoning vs. Documentation The Next Battleground for Medical LLMs

MAR 20, 2026 • PAID



Abstract

The first wave of healthcare AI scored decisive wins in documentation automation. Ambient scribes, coding copilots, and summarization layers delivered clear ROI solving a well-bounded problem: compress high-entropy clinical inputs into structured, billable outputs. That layer is now saturating. The next frontier is harder, more valuable, and genuinely unsolved: augmenting clinical reasoning itself.

This essay covers:

- Why documentation AI succeeded and why reasoning AI is fundamentally different
- The three hard architectural requirements current LLMs only partially meet (structured representation, hypothesis generation, uncertainty quantification)
- Why next-token predictors structurally struggle with clinical cognition
- Emerging architectures trying to bridge the gap (tool-augmented reasoning, graph-based inference, persistent memory layers)
- Failure modes unique to reasoning-adjacent systems
- Why current benchmarks like MedQA are nearly useless for evaluating actual reasoning
- The economic argument for why reasoning AI creates durable moats that documentation AI cannot

- A framework for thinking about AI's role: Advisor vs. Cognitive Extender vs. Autonomous Reasoner

Table of Contents

The Documentation Win and Why It's Running Out

Compression vs. Inference: A Real Distinction

Three Requirements That Break Current LLMs

Why the Architecture Itself Is the Problem

What's Actually Being Built to Fix This

The Failure Modes Nobody Talks About

Evaluation Is Broken and Everyone Knows It

Three Paradigms for AI's Role in Reasoning

The Economic Case for Betting on Reasoning

What a Computable Differential Actually Looks Like

The Documentation Win and Why It's Running Out

If you've been paying attention to where healthcare AI dollars have gone over the last four years, the pattern is pretty obvious. Ambient scribes, prior auth automation, clinical note summarization, revenue cycle coding assist. Every major health system has piloted at least one of these. Most are in some phase of deployment. Nuance, Abridge, Suki, Nabla, and a handful of EHR-native products from Epic and Oracle have collectively reshaped how clinicians think about administrative burden. That's not nothing. It's actually a big deal.

The ROI story for this category is clean and defensible. Physicians were spending somewhere between one and three hours per day on documentation depending on specialty. Ambient scribes demonstrably cut that. KLAS Research data from 2022 showed DAX users saving an average of 7 minutes per note. Multiply that across patient day and you're talking real productivity gains. Payers and health systems quantify it. CFOs could model it. Procurement decisions got made.

But here's the uncomfortable reality underneath that success: documentation AI is fundamentally a compression problem. It takes high-entropy inputs, which is to say the rambling, overlapping, sometimes contradictory content of a clinical encounter and transforms them into low-entropy structured outputs. Progress notes. HCC codes. Discharge summaries. The model doesn't need to understand what's clinically happening. It needs to recognize patterns in language and map them onto the expected structure of a clinical document. That's a solved class of problem. Transformers are exceptionally good at it.

The saturation dynamic is already visible in the market. Ambient scribe functionality is becoming table stakes. Epic and Oracle are bundling it natively. Gross margins are compressing. The differentiation thesis for pure-play documentation vendors is getting harder to sustain. That doesn't mean the category is dead. It means the wave has crested and the smart capital is looking at what comes next.

Compression vs. Inference: A Real Distinction

The phrase "clinical reasoning" gets thrown around a lot, sometimes loosely, so it's worth being precise about what it actually means and why it's a different beast from documentation.

Documentation is a compression problem. Clinical reasoning is an inference problem. Those are not the same thing, and conflating them has led to a lot of overpromising in the AI health space.

Here's what inference under clinical uncertainty actually involves. A 58-year-old chest pain and new dyspnea walks into the ED. The clinician is not searching a knowledge base. They are constructing a probabilistic model in real time. What's the prior probability of ACS in this demographic given this symptom cluster? How does the troponin trend update that probability? What does the absence of pleuritic component tell me about PE likelihood? If the D-dimer comes back mildly elevated in the context of a recent long flight, how much does that shift things? This is Bayesian updating applied under time pressure across a partially observed, dynamically evolving data set.

Documentation systems don't need to do any of this. They just need to accurately capture and reformat what happened. The clinician already did the reasoning. The scribe just records it.

The reason this distinction matters for investors and builders is that it defines the time and durability of the value being created. Compressing a clinical note saves time. Improving diagnostic inference changes outcomes. And changing outcomes is what healthcare actually spends money on. The US misdiagnosis rate hovers around 12 million adults per year according to data published in BMJ Quality and Safety. Diagnostic error contributes to somewhere between 40,000 and 80,000 deaths annually by most estimates. The economic footprint of that problem dwarfs the administrative burden of the problem by a considerable margin.

Three Requirements That Break Current LLMs



Continue reading this post for free, courtesy of Special Interest Media.

[Claim my free post](#)

Or purchase a paid subscription.

← Previous

Next

© 2026 Thoughts on Healthcare · [Privacy](#) · [Terms](#) · [Collection notice](#)
[Substack](#) is the home for great culture