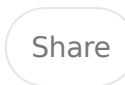
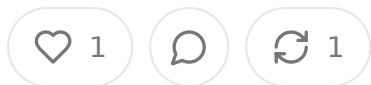


The Pattern Always Repeats: Why Healthcare's Next Revolution Runs on Electricity, Not Software

MAR 02, 2026 • PAID



Abstract

This piece argues that the three horsemen of every major economic revolution have always been communication, energy, and transportation, and that understanding historical sequencing helps predict where healthcare goes from here.

Key claims:

- Every major economic upheaval from the printing press forward follows the same three-part unlock pattern
- LLMs, specifically GPT-4 and successors, represent healthcare's communication unlock, equivalent in magnitude to Gutenberg or Morse
- Healthcare has had its communication revolution. The energy revolution is coming next and will dwarf what software alone can do
- Nvidia has quietly become one of the most important energy infrastructure companies on the planet, and almost nobody is framing it that way
- Quantum computing, nuclear fusion, next-gen electrical infrastructure, and AI efficiency gains are the mechanisms
- For investors and founders, the implication is that the 2030s will look nothing like the 2020s in terms of what healthcare can actually compute, model, and deliver in time

The time to position is before the energy unlock, not after

Table of Contents

The Three-Part Pattern (and why most people miss it)

The Printing Press Didn't Save Lives, But It Started the Chain

Steam, Coal, and the First Time We Industrialized Medicine

Electricity, Railroads, and the Birth of the Modern Hospital

The Internet as a Communication Unlock and Why Healthcare Barely Felt It

LLMs Are Healthcare's Gutenberg Moment

Why Communication Alone Never Finishes the Job

Nvidia and the Quiet Energy Revolution Inside the Chip

The Energy Unlock Is the Missing Piece

What Quantum and Fusion Actually Mean for Healthcare (Practically)

How to Invest Ahead of an Energy Revolution You Can't Fully Predict

The Three-Part Pattern (and why most people miss it)

Historians love to argue about what causes economic revolutions. Was it a charismatic leader? A lucky war outcome? A policy shift? Occasionally it was all of those things, but underneath virtually every transformational economic leap in the last five centuries, you find the same boring trio doing the heavy lifting: something changed in how people communicated, something changed in how they generated or moved energy, and something changed in how they moved physical things through space. Communication, energy, transportation. Every time. Like clockwork, except the

runs on about a hundred-year cycle, which is inconveniently longer than most investment horizons.

The reason most people miss the pattern is that they tend to fixate on the sexy individual invention, the printing press, the steam engine, the microchip, and the internet like a standalone miracle. But none of those things worked in isolation. The printing press mattered because it happened alongside early capitalism and paper supply chains. The steam engine mattered because coal extraction had gotten good enough to actually fuel it consistently. The internet mattered because fiber optics, server farms, and the end of the Cold War all conspired to make it globally deployable. When you zoom out far enough, the individual invention looks less like a cause and more like a symptom of three underlying systems converging at once.

Healthcare is not exempt from this pattern. In fact healthcare is one of the best case studies available for understanding how the pattern plays out in a domain that is simultaneously information-intensive, energy-hungry, and physically distributed. When you look at the history of medicine through the lens of communication, energy, and transportation, you get a completely different and arguably more useful story than the standard “great scientists and great discoveries” narrative.

The Printing Press Didn't Save Lives, But It Started the Chain

Johannes Gutenberg finished his press somewhere around 1440, give or take some academic debate, and it would be historically generous to say it immediately improved healthcare outcomes. It did not. What it did was democratize the written word, which meant that medical texts, previously copied by monks at enormous expense and distributed to maybe a few dozen institutions across Europe, could now be printed in the hundreds and eventually thousands. Andreas Vesalius published “De Humana Corporis Fabrica” in 1543 with detailed anatomical illustrations and it spread across European universities in a way that would have been physically impossible a century earlier.

This is the communication unlock doing its first pass on medicine. The knowledge existed in fragments before Gutenberg. What printing did was aggregate, standardize, and distribute it. Physicians in Padua and physicians in Paris could now argue from the same text, which sounds trivial until you realize that standardized knowledge sharing is the prerequisite for anything resembling evidence-based practice. You cannot have a scientific consensus if nobody can read each other's work.

But notice what printing could not do. It could not make surgery safer because it could not sterilize instruments. It could not speed up the transportation of sick patients. It could not power the equipment that did not yet exist. The communication unlock planted seeds that would take centuries to fully germinate because the energy and transportation unlocks had not yet arrived. The lag between a communication revolution and its full downstream effects on something as complex as healthcare is not a bug, it is a feature of how these systems build on each other.

Steam, Coal, and the First Time We Industrialized Medicine

The first industrial revolution, roughly 1760 to 1840, was predominantly an energy story. Coal became cheap and extractable at scale, the steam engine became reliable enough to run factories, and for the first time in human history, productive output decoupled from human or animal muscle. The economic consequences were staggering and well-documented. The healthcare consequences were, for a while, genuinely terrible.

Urbanization happened faster than sanitation infrastructure could accommodate. Workers flooded into cities that had no sewage systems, no clean water, and no understanding of germ theory. Life expectancy in English industrial cities during the early 1800s was lower than in rural areas, which is a remarkable historical fact given that cities had more access to doctors and hospitals. The energy revolution created conditions for epidemic disease faster than medicine could respond to it.

But then the energy revolution started to pay its forward-looking dividends. Rail networks, powered by steam, meant that pharmaceutical compounds and medical

supplies could move across national territories in days rather than weeks. The second industrial revolution, coal-powered industrialization that was killing urban workers was also enabling the mass manufacturing of medical instruments at a scale and quality level that artisan craftsmen could not match. Surgical tools got better and cheaper. Hospitals could now source supplies reliably. The first real pharmaceutical supply chains emerged.

The transportation unlock arrived right on schedule, hand in hand with the energy unlock, and together they started converting the communication unlock of the printing press into something operationally useful. The diffusion of medical knowledge that Gutenberg made possible started to translate into standardized clinical practice because hospitals could now get the same instruments, the same compounds, and train physicians from the same texts. The three systems were finally talking to each other.

Electricity, Railroads, and the Birth of the Modern Hospital

The second industrial revolution, from roughly 1870 to 1914, is where healthcare really started to look like something modern. This one had a different energy profile: electricity rather than steam, and the internal combustion engine rather than railroads as the primary transportation unlock. The communication unlock of this era was the telegraph and then the telephone, which meant that for the first time you could transmit information at something close to the speed of thought rather than the speed of horses.

The hospital as an institution we would recognize today is fundamentally a product of this era. Electric lighting made surgery viable after dark and in poorly lit spaces. Autoclaves powered by electric or steam heat made sterilization systematic rather than occasional. The telephone meant that physicians could consult with specialists in real time across distances that would have previously required a multi-day journey. X-rays, discovered in 1895, required reliable electrical current to operate. The EKG, developed in the early 1900s, was entirely dependent on electrical infrastructure.

It is worth pausing on how much of what we take for granted in clinical medicine is essentially an artifact of the second electrical revolution rather than a product of scientific genius. The germ theory of disease was understood intellectually before Pasteur and Koch formalized it, but it became clinically actionable only when hospitals had the energy infrastructure to actually sterilize things reliably. The knowledge arrived before the energy capability. The energy capability converted knowledge into practice. This is the pattern doing exactly what the pattern always does.

Transportation in this era deserves its own moment because the emergence of rail and early automobile networks changed the geography of healthcare delivery in ways that are still structurally embedded in the system today. The hospital as a centralized institution only makes sense if patients can travel to it reliably. Before reliable transportation, healthcare was necessarily distributed because you could move sick people safely at any scale. Rail and automobile networks made centralization economical, and the modern hospital system was built on that assumption. Spoiler: that assumption is about to get challenged again.

The Internet as a Communication Unlock and Why Healthcare Barely Felt It

The internet represents what should have been healthcare's third major communication unlock, following the printing press and the telegraph, and by many measures it underperformed its potential spectacularly. The commercial internet arrived in the early 1990s. By the mid-2000s, almost every other information-intensive industry had been structurally transformed. Music, retail, financial services, travel, and media were all fundamentally reorganized around digital information flow. Healthcare spent twenty years fighting over fax machines.

The reasons for healthcare's internet-era lag are well-rehearsed among the investors and operators reading this, so no need to dwell on them extensively. Regulatory structure, reimbursement misalignment, liability concerns, the sheer complexity of existing workflows, all contributed. But there is a structural reason that goes beyond

the policy and economic friction. The internet as a communication technology was built for moving structured data and documents efficiently. Healthcare's core problem has never been moving structured data efficiently. Its core problem has been interpreting unstructured clinical information, the physician note, the pathology report, the radiology read, the patient narrative, and doing something clinically with it. EHRs got good at storing documents. They got terrible at understanding them.

What digital health as an industry did from roughly 2010 to 2022 was largely take existing clinical workflows, digitize them, and add a consumer-facing layer. That's not nothing, but it is a long way from the kind of structural transformation that happened to music or financial services. The communication unlock the internet provided was a translation of analog documents into digital ones, which improved efficiency at the margins but did not change the fundamental epistemological problem of clinical medicine, which is that most of the meaningful information in healthcare is embedded in language and context that requires genuine interpretation rather than just retrieval.

LLMs Are Healthcare's Gutenberg Moment

GPT-4 launched in March 2023. Within six months, every serious health tech investor had at least ten portfolio companies claiming they were building the AI layer for healthcare. Within a year, the more thoughtful founders had figured out that the interesting question was not whether LLMs could pass the USMLE (they could, but rather what specific clinical tasks they could automate reliably enough to truly scale. By 2024 and into 2025, the actual use cases started separating into things that clearly worked and things that clearly did not.

What clearly worked: prior authorization automation, clinical documentation reduction, patient-facing triage and navigation, coding and revenue cycle optimization, literature synthesis for clinical decision support. What remained but anything requiring genuine reasoning under uncertainty with life-or-death stakes

anything involving subtle physical examination findings that have no digital representation, anything where model hallucination was clinically intolerable. That pattern holds. LLMs are a communication unlock, not a universal problem solver.

Here is the historical analogy that is most useful. The printing press did not immediately cure any diseases. What it did was create the infrastructure for scientific consensus-building that eventually made evidence-based medicine possible about 100 years later. LLMs are doing something structurally similar on a compressed timeline. They are not curing diseases in 2024. What they are doing is creating the infrastructure for a completely different relationship between clinical knowledge and clinical action. For the first time, unstructured clinical language can be interpreted, synthesized, and acted upon at computational speed. That is the printing press moment for medicine, not the cure moment.

The investor implication of framing it this way is important. Companies that are building on LLMs right now are building during the Gutenberg era, which was genuinely transformative but which also preceded the steam engine by a couple centuries. The returns from communication-layer companies will be real and some will be generational. But the companies that capture the majority of the economic value from this broader revolution will likely be companies that combine the communication unlock with the next energy unlock, and we are just now starting to see what that looks like.

Why Communication Alone Never Finishes the Job

There is a useful exercise here which is to ask what the specific computational and physical limitations are that prevent the full deployment of what LLMs can theoretically do in healthcare. Not the regulatory limitations, those will resolve over time, and not the workflow limitations, those will also resolve. The genuine physical limitations.

The first one is compute cost. Running a frontier model at clinical scale, meaning continuous inference across a large health system's entire patient population in :

time, is currently prohibitively expensive. The math on GPU inference costs versus reimbursement rates in most clinical workflows does not close unless you are in high-value specialty or a revenue cycle application. This is not a software problem, it is an energy and hardware problem. The cost of inference will fall as hardware improves, but the rate of improvement is constrained by the energy density of current semiconductor architectures.

The second one is data latency. Clinical AI that actually improves outcomes in acute care settings needs to process information faster than current infrastructure allows. An ICU early warning system that processes vitals every fifteen minutes is better than nothing, but genuinely transformative clinical AI needs continuous real-time processing of physiological signals, imaging, lab values, and contextual clinical data simultaneously. The compute demands for this, at scale, across an entire health system, are currently not achievable with economically viable infrastructure.

The third one is the physical environment. A significant fraction of the highest-value clinical applications for AI involve scenarios where the AI needs to be present at the point of care in a physically demanding environment, the operating room, the ambulance, the remote clinic, and where current power infrastructure is either unavailable or insufficient. Inference at the edge requires local compute. Local compute requires power. Power in clinical edge environments is a serious engineering constraint that is not being solved by software.

All three of these limitations point at the same underlying bottleneck: energy density and energy cost. This is not a coincidence. It is the pattern telling you something

Nvidia and the Quiet Energy Revolution Inside the Chip

Before getting to fusion and quantum, there is a company already mid-revolution that does not get nearly enough credit for being an energy story rather than just a chip story. Nvidia is, depending on how you want to think about it, either the most important infrastructure company of the current decade or the most underappreciated energy efficiency company in history. Probably both.

The H100, Nvidia's flagship data center GPU released in 2022, delivers roughly 200 teraflops of FP16 performance at around 700 watts. That sounds like a lot of power until you compare it to what the same compute workload would have required five years earlier. The A100, released in 2020, was already a massive leap over its predecessors. The H100 was another leap over that. The trajectory of compute-per-watt across Nvidia's data center GPU generations over the last decade is roughly doubling every two to three years, which is actually faster than classical Moore's Law and represents one of the most significant energy efficiency improvements in the history of computing.

The Blackwell architecture, released in 2024, pushed this further with the GB200 NVLink system delivering roughly five times the inference performance of H100 in roughly comparable power envelopes for many workloads. For large language model inference specifically, Nvidia has published benchmarks showing the GB200 delivering around 30 times better performance per watt on certain inference tasks compared to the H100 generation. Even discounting for benchmark optimism, the directional story is real and it is large.

Why does this matter specifically for healthcare AI? Because the economics of clinical AI deployment are directly downstream of inference cost, and inference cost is directly downstream of compute-per-watt. Every time Nvidia's efficiency curve moves, use cases that were previously economically unviable become viable. The continuous ICU monitoring application that did not pencil out in 2022 might pencil out in 2025 on Blackwell infrastructure. The real-time surgical assistance system that requires a cost structure unavailable outside of major academic medical centers might be deployable at community hospital scale within the next two GPU generations.

There is also a subtler point worth making about what Nvidia has done for the energy consumption profile of AI at the national level. Data center electricity consumption in the US has been growing rapidly alongside AI workload growth, and projections from groups like the Lawrence Berkeley National Laboratory have estimated that US data centers could consume somewhere between 6 and 12 percent of national electricity by the late 2020s depending on efficiency improvement trajectories. Nvidia's efficiency gains are one of the primary mechanisms by which that growth is being partially

offset. The company is simultaneously making AI more powerful and making AI energy-intensive per unit of output, which is the exact combination of factors you need to make clinical AI at scale economically rational.

Nvidia has also made less-publicized but important investments in the software for energy-efficient inference. The TensorRT inference optimization library, the microservices architecture announced in 2024, and the specific optimizations Nvidia has built for healthcare AI workloads through its Clara platform are all part of a coherent strategy to make the GPU do more clinical work per kilowatt-hour. Jen Huang has been explicit in public forums that Nvidia views itself as an energy efficiency company as much as a semiconductor company, and the product roadmap supports that framing more than the typical chip company narrative does.

The healthcare-specific implication of Nvidia's trajectory is worth isolating. Nvidia's Clara Parabricks platform, which handles genomic sequencing computation, reduced the time to run a full genome variant analysis from roughly 30 hours on CPU to 30 minutes on GPU, at dramatically lower energy cost per analysis. At current sequencing volumes that is a meaningful efficiency gain. At the sequencing volumes projected for the late 2020s as genomic medicine becomes more mainstream, the difference between CPU and GPU infrastructure for genomics is the difference between the economics working and not working. Nvidia is not incidentally part of that story. It is central to it.

None of this means Nvidia is risk-free as a bet or that the current valuation is obviously rational. It means that when thinking about the energy unlock for healthcare AI, Nvidia belongs in the same conversation as fusion and quantum, because it is already delivering energy efficiency improvements that are changing what is computationally possible in clinical settings right now, not in fifteen years.

The Energy Unlock Is the Missing Piece

There are three additional energy technologies in active development that have genuine potential to change the constraints described above within the next decade or fifteen years: next-generation nuclear fission via SMRs, nuclear fusion, and quantum computing.

computing as a distinct computational energy paradigm. Together with Nvidia's efficiency curve, these represent four different mechanisms converging on the same outcome: dramatically cheaper and more abundant compute for clinical AI.

SMRs are the nearest-term story on the generation side. Companies like NuScale, Kairos Power, and a handful of others are building reactors designed to deploy at a small scale with dramatically reduced construction timelines compared to traditional nuclear. The DOE has been funding these programs aggressively and several are in late-stage regulatory review. For healthcare specifically, the implication of SMR deployment is not directly about powering hospitals. It is about powering the data center infrastructure that clinical AI depends on. If the marginal cost of electricity for large-scale computing falls significantly, the economics of clinical AI at scale change dramatically, and they compound with the efficiency gains Nvidia is already delivering on the demand side.

Fusion is a longer bet. Commonwealth Fusion Systems, with its SPARC project, achieved a significant magnet milestone in 2021 and has credibly moved the expected timeline for net energy gain to sometime in the late 2020s or early 2030s. TAE Technologies, Helion Energy backed by Sam Altman (make of that what you will the LLM connection), and a handful of others are pursuing different confinement approaches. The probability of commercially viable fusion electricity by 2035 is no longer negligible. If it arrives, the downstream effects on compute-intensive applications like clinical AI are hard to overstate, and they land on top of an Nvidia efficiency curve that will have continued improving in the interim.

Quantum computing is the most speculative of the three but deserves inclusion because its implications for healthcare go beyond cheaper computation. Quantum systems have theoretical advantages in molecular simulation that classical computers will never match efficiently. Drug discovery, protein folding beyond what AlphaFold can currently do, and the optimization problems involved in personalized dosing and treatment sequencing are all domains where quantum advantage, when it arrives, could compress development timelines by years. The energy story with quantum is slightly different since current quantum systems require extreme cooling to near

absolute zero which is itself an enormous energy cost, but that constraint is exactly what the next generation of room-temperature qubit approaches is designed to solve.

The healthcare energy unlock is not a single technology. It is a cluster of converging developments in how we generate, transmit, and apply energy to computation and physical environments, arriving more or less simultaneously within a fifteen-year window. Nvidia is the part of that cluster that is already here and already changing the math. SMRs, fusion, and quantum are the parts that are loading. This is exactly like the second industrial revolution worked. Electricity, the internal combustion engine, and the telephone were not planned to converge. They just did, because the underlying physics and engineering were all mature enough at the same historic moment.

What Quantum and Fusion Actually Mean for Healthcare (Practically)

Getting concrete matters because the history lesson is useful but ultimately investors and founders need to know what to actually build or fund.

The first practical implication is that the cost curve for clinical AI inference will be much faster than most current DCF models assume. Healthcare AI company valuations based on current GPU costs as a structural assumption are going to look wrong within a few years, probably in a positive direction for pure AI application companies and in a complicated direction for companies whose moat is primarily compute efficiency at current energy prices. Nvidia's roadmap alone would drive this conclusion. Add quantum deployment to the grid and it accelerates.

The second practical implication is that real-time physiological AI becomes economically viable in a low-energy-cost world. Right now the most interesting continuous monitoring applications, the ones that can genuinely predict a sepsis or a cardiac deterioration hours before current early warning systems, are not deployed at scale because the compute cost of continuous inference on rich sensor data does not pencil out against DRG reimbursement. Drop the inference cost by eighty percent, which is a plausible outcome of the combined Nvidia efficiency c

plus grid energy improvements over the next decade, and that math changes. The companies building sensor infrastructure and clinical workflow integration for IoT applications right now are building something that will be worth dramatically more when the energy unlock arrives.

The third practical implication is drug discovery timelines. Quantum molecular simulation, even at relatively early fidelity levels, will let computational chemists run experiments in silico that currently require months of wet lab work. The companies building platform chemistry capabilities with serious computational infrastructure assumptions are positioning for a world where the rate-limiting step in drug development shifts from synthesis and testing to regulatory and clinical validation. That is a fundamentally different R and D economics problem with different capital requirements and different risk profiles.

The fourth implication, and this one is underappreciated, is what cheap abundant energy means for healthcare delivery infrastructure in the developing world. A significant fraction of the global disease burden sits in environments where healthcare delivery is constrained not just by communication and knowledge but by the physical impossibility of running diagnostic equipment reliably. Point-of-care diagnostics require consistent power. Cold chains for vaccines and biologics. Hospital-level surgical capability in rural settings. These are not software problems. They are energy problems. A world with dramatically cheaper and more distributed energy production is a world where a lot of global health infrastructure problems become tractable in ways they simply are not today. Notably, Nvidia's edge inference platforms are part of this story too: the ability to run meaningful clinical AI models on low-power edge hardware in resource-constrained environments is already improving and will continue to.

How to Invest Ahead of an Energy Revolution You Can't Fully Predict

The honest version of this section starts with admitting that the timing of the full energy unlock is genuinely uncertain. Fusion could be 2032 or it could be 2045.

regulatory timelines have a way of extending. Quantum decoherence problems may take another decade to crack at room temperature. What is not uncertain is Nvidia's efficiency trajectory, which is documented, roadmapped, and already changing the economics of clinical AI deployment in real time. That is the part of the energy story you can bet on now without requiring a fusion breakthrough.

Given that range of uncertainty across the different mechanisms, the sensible investment posture is to look for healthcare companies whose value proposition becomes dramatically stronger in a world of cheap abundant compute and cheap abundant energy, and to prefer those companies when they can generate returns on the slower timeline. Find the businesses where the energy unlock is option value rather than the base case.

Clinical AI application companies in revenue cycle, prior auth, and documentation already generating real returns on current energy economics. They get meaningfully better with cheaper inference. That is a good risk profile and it improves with every Nvidia generation regardless of what happens with fusion.

Continuous monitoring and early warning companies with serious data and world integration built already are interesting asymmetric bets. The underlying clinical value is not in question. The economics are. An energy unlock, whether it comes primarily from Nvidia's efficiency curve or from grid-level changes, resolves the economics.

Computational drug discovery platforms that have built serious proprietary chemical data and molecular modeling infrastructure are positioning correctly for quantum advantage. The ones generating near-term revenue from classical compute applications on current Nvidia infrastructure have a survivable intermediate timeline.

Healthcare infrastructure companies focused on emerging markets, especially those building modular or low-power diagnostic and delivery capabilities, are playing a game that the energy unlock potentially pulls forward by a decade.

What to be cautious about: healthcare AI companies whose primary defensibility claim is inference efficiency or proprietary compute optimization. If Nvidia's

efficiency curve plus eventual grid improvements commoditize inference, that market narrows or disappears. Also worth scrutinizing: any business model that requires current energy constraints to persist in order to maintain competitive advantage.

The meta-lesson from the historical pattern is that the people who made the most progress from the second industrial revolution were not primarily the people who invented new technologies. They were the people who built the applications, systems, and institutions designed for a world powered by electricity before most of the world was. The hospital system, the pharmaceutical supply chain, the medical device industry—all of these were built by people who assumed electricity was real and was coming and designed forward from that assumption.

Nvidia is telling you, with published silicon roadmaps and real benchmark data, the next several years of compute economics look like. Fusion and quantum are telling you, with less certainty but with increasing credibility, what the decade after that might look like. The same opportunity exists right now for the generation of healthcare tech founders and investors willing to design forward from the assumption that energy unlock is real and is compounding. The ones who assume a world of dramatically cheaper and more abundant compute, who assume that the physical constraints on clinical AI will be lifted within their company's scaling timeline, who assume that the geographies currently excluded from modern healthcare by energy constraints will be included: those are the founders building the institutions that will define healthcare delivery in the second half of the 21st century.

The pattern always repeats. The communication unlock just landed. The energy unlock is already partly here and the rest of it is loading.



1 Like • 1 Restack

[← Previous](#)

[Next](#)

Discussion about this post

Comments

Restacks



Write a comment...

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