# AI First Principles Treatise

Designing Organizational Artificial General Intelligence

## The Case for the Principles

This treatise is the definitive collection of evidence and reasoning that underpins the AI First Principles. It serves as the formal, foundational proof for the framework, not a guide for its implementation.

## Introduction

The discourse surrounding artificial intelligence is dominated by a grand and misleading narrative of replacement. We are told AI is coming for jobs, that human obsolescence is a matter of when, not if. This is a failure of imagination. The true, immediate potential of AI is not replacing people, but replacing the soul-crushing bureaucracy that makes their work miserable. It offers a path to dismantle the layers of process, approvals, and communication overhead that drain human potential.

Yet, this opportunity is being squandered. A staggering percentage of AI projects stall before deployment, with studies showing that over 80% fail to move from the lab to production.[1] The reasons are rarely technical. The failure is human. Organizations are implementing AI backward, attempting to automate inefficiency by bolting a revolutionary capability onto broken operational models. They are strapping a jet engine to a horse-drawn carriage. The result is not a faster carriage; it is a catastrophic mess, a predictable wave of costly implementations that fail to deliver value.

This paper introduces the AI First Principles, a framework for reversing this approach. These are not another methodology, but a collection of foundational principles for rethinking work in an organization where AI is a native collaborator. These principles are grounded in decades of research from disparate fields: the systems thinking of Peter Senge,[2] the human-centered design philosophy of Don Norman,[3] the practical operationalization of AI detailed in *Age of Invisible Machines*,[4] the operational logic of Eli Goldratt's Theory of Constraints,[5] and the innovation theory of Clayton Christensen.[6]

A critical insight emerged from conversations with Dave Thomas, co-author of the Agile Manifesto: the principles must reveal consequences, not prescribe solutions. As Thomas noted during his review, "AI behaves more like people than engineering"—it's unpredictable, context-dependent, and prone to the same biases that plague human systems.[7] This treatise therefore focuses on what breaks when these principles are ignored, not on what to do. Prescription dates quickly; understanding consequences creates timeless wisdom.

The framework presented here offers a coherent, human-centered approach for any organization ready to stop automating the past and start building the future of work.

# Beyond These Principles

The AI First Principles are intentionally focused on the timeless truths of human-AI collaboration. They are not a comprehensive guide to every facet of implementation. These principles operate on the assumption that specialized domains like regulatory frameworks, legal liability, and specific technology choices are already being handled by their respective experts. For example, the NIST AI Risk Management Framework provides a robust, voluntary structure for identifying, measuring, and monitoring risks, serving as an excellent operational counterpart to the principles outlined here.[8] These principles are designed to be layered upon this foundation of responsible governance, not to replace it. They are for the builders, designers, and strategists tasked with creating value within those established boundaries, ensuring that our most powerful tools are fused with human ingenuity to eliminate dysfunction and unlock potential.

# The Principles

### ***AI Inherits Messiness***

AI learns from people. Therefore AI systems, like people, are inconsistent and operate more effectively with structure. Trying to engineer them to operate like deterministic code will result in system failures. Variation is inevitable, not accidental. *Define what's prohibited over what's required.*

> **The Hidden Problem** The most dangerous assumption in AI implementation is that algorithms are objective. This belief creates a false sense of security that allows bias to scale invisibly. When an AI system learns from human-generated data—hiring decisions, loan approvals, medical diagnoses—it doesn't extract pure truth. It extracts patterns, including the prejudices, shortcuts, and contextual assumptions embedded in those decisions. The hidden problem is that these inherited patterns operate with machine speed and scale, turning what might have been isolated human errors into systematic discrimination. This is not a technical flaw to be patched; it is a fundamental characteristic of how machine learning works. As Cathy O'Neil documents in *Weapons of Math Destruction*, when biased data trains predictive models, those models don't just reflect existing inequities—they amplify and legitimize them under the guise of algorithmic neutrality.[^9]

> **The Cost of Ignoring It** When organizations fail to recognize that AI inherits human patterns, they build systems that perpetuate historical discrimination at unprecedented scale. Amazon's recruiting AI, trained on a decade of predominantly male hiring decisions, learned to penalize resumes that included the word "women's"—a direct translation of past bias into algorithmic policy.[^10] The cost is not just legal liability or reputational damage; it's the systematic exclusion of talent and the erosion of trust in automated systems. More insidiously, the mathematical veneer of AI makes these biases harder to challenge. A human manager's discriminatory decision can be

questioned; an algorithm's output is often treated as fact. This creates what Safiya Noble calls "algorithmic oppression"—where technology encodes and enforces social inequalities while appearing neutral and inevitable.[^11] Organizations that ignore this principle don't just fail ethically; they build brittle systems that collapse when their inherited assumptions no longer match reality.

**The Contrarian Insight** The conventional approach to AI bias focuses on cleaning the data and adjusting the model—technical fixes for what is fundamentally a design problem. The contrarian insight is that variation in AI output is not a bug; it's a feature of learning from messy human reality. The goal should not be to eliminate all variation (which is impossible) but to establish explicit constraints that prevent harmful patterns from manifesting. This shifts the focus from achieving perfect objectivity—which is unattainable—to defining acceptable boundaries. As Timnit Gebru and others argue in their work on datasheets for datasets, the solution is not sanitizing training data but understanding and documenting what biases it contains, then designing systems with guardrails that prevent those biases from causing harm.[^12] The insight is that constraint is not a limitation on AI's capability; it's the prerequisite for its responsible use.

**Solution Framework** Implementing this principle requires acknowledging that every AI system inherits the context of its training data and then deliberately designing constraints around that inheritance. This means starting with a "bias inventory"—a systematic audit of what patterns the training data likely contains. Who created this data? Under what circumstances? What was considered normal or successful in that context? Once these inherited patterns are visible, the next step is to establish operational boundaries: rules, thresholds, and review triggers that prevent harmful outputs from reaching production. This is not about achieving algorithmic perfection but about building "circuit breakers" into the system. For example, if a hiring AI shows statistically significant demographic disparities in its recommendations, the system should flag those patterns for human review rather than auto-executing them. This approach mirrors the principle of "bias bounties" proposed by AI ethics researchers, where organizations actively seek to discover and document failure modes before they cause harm.[^13]

**Philosophical Foundation** This principle is rooted in the philosophy of situated knowledge, which argues that all knowledge—including data—is produced from a particular perspective and carries the assumptions of its context.[^14] Algorithms trained on that data are not discovering universal truths; they are learning the patterns of a specific time, place, and power structure. This connects to the concept of "algorithmic accountability," which holds that the designers and deployers of AI systems must acknowledge and take responsibility for the value judgments embedded in those systems.[^15] Furthermore, the inevitability of variation aligns with the Bayesian understanding of uncertainty in statistical inference: all predictions carry inherent uncertainty that should be quantified and managed, not hidden.[^16] By

acknowledging that AI inherits human patterns, we commit to building systems that make those inheritances visible and manageable rather than pretending they don't exist.

Implementation Reality The primary resistance to this principle comes from the desire for algorithmic determinism—the comforting but false belief that machines can provide objective answers. Teams will push back against constraint mechanisms, arguing they reduce efficiency or limit AI's potential. This resistance must be met with evidence: show the catastrophic failures that occur when inherited bias scales unchecked. The practical implementation requires building bias detection into the development pipeline as a mandatory gate, not an optional review. This means establishing clear metrics for acceptable variation (e.g., demographic parity in outcomes) and treating violations of those metrics as seriously as security vulnerabilities. It also requires a cultural shift: treating variation as data rather than failure, and building feedback systems that allow the organization to learn from the patterns the AI reveals about its own historical decisions.

# ***AI Fails Silently***

AI accumulates errors across thousands of interactions before patterns become visible. Traditional systems failed loudly with clear signals; AI fails quietly on repeat. _Build feedback loops, not post-mortems._

The Operational Dysfunction The modern technology ethos is dominated by the mantra of "move fast and break things." This iterative, fail-fast approach, popularized by the Agile movement, is incredibly powerful for navigating uncertainty in product features and user interfaces. However, when misapplied to AI systems operating at scale, it becomes a recipe for disaster. Traditional systems failed in ways humans could observe and intervene: a clerk processing a hundred forms would catch anomalies, a manager reviewing daily reports would spot trends. AI systems process thousands of decisions per second, compounding errors before any human realizes something is wrong. You cannot "iterate" your way to safety after an AI has already denied ten thousand legitimate insurance claims based on a flawed assumption. You cannot "A/B test" your way out of a trading algorithm that has crashed a market. The operational mistake is treating AI-scale risks as if they are the same as human-scale mistakes—assuming that the feedback loop that works for human processes will work when those processes run at machine speed.

Why Standard Approaches Fail Standard software development lifecycles, even agile ones, are built on the assumption of observable failure. When something breaks, someone notices, and the team can respond. This assumption crumbles when AI operates across thousands of edge cases simultaneously, many of which fall outside the training scenarios. An AI trained to approve loans might encounter a novel economic condition and systematically misclassify an entire demographic group—not

because it's biased against them, but because it has no reference for how to evaluate risk in the new context. By the time the pattern becomes visible in aggregate statistics, the damage is done. The standard approach of "deploy and monitor" fails because monitoring operates on human timescales while failure propagates at machine speed. This mirrors the insight from high-reliability organizations (HROs) like nuclear power plants: when the cost of failure is catastrophic, you cannot rely on learning from mistakes; you must engineer the system to be fundamentally safe by design.[^27]

**The Principle's Core Logic** The core logic of this principle is to acknowledge a fundamental asymmetry: AI can cross boundaries—ethical, operational, legal—thousands of times in the time it takes a human to recognize there's a boundary to cross. This doesn't mean AI is inherently dangerous; it means ambiguous authority becomes catastrophic delegation at machine speed. When you hand a human employee vague instructions, they'll ask clarifying questions. When you hand an AI system vague instructions, it will execute its interpretation millions of times before you notice the interpretation was wrong. The solution is not to slow AI down to human speed—that defeats its purpose—but to invest heavily upfront in defining boundaries and validating capability before granting authority. This means treating AI deployment not as iteration but as a controlled release: start with narrow, well-defined domains where you've validated the AI's behavior, then expand deliberately only after proving it handles edge cases correctly. As Nassim Taleb argues in *Antifragile*, some systems benefit from volatility and error (antifragile), while others are destroyed by it (fragile). AI systems operating at scale are inherently fragile to unanticipated inputs—they don't learn from catastrophic errors the way humans do; they simply execute them faster. [^38]

**Research Validation** This principle is grounded in research on human-automation interaction and specifically the concept of "automation surprise"—when an automated system behaves in ways operators don't expect, often with catastrophic results.[^39] Aviation research has documented numerous cases where autopilot systems exceeded their design parameters faster than pilots could intervene. The same dynamic applies to AI: the speed at which decisions compound means that by the time humans notice something is wrong, the system may have already caused irreversible harm. Research on algorithmic accountability emphasizes the need for "circuit breakers"—automated mechanisms that halt AI operation when it detects it's operating outside validated parameters.[^40] This principle also draws from the DevOps concept of "chaos engineering," where systems are deliberately stressed to discover failure modes before they manifest in production,[^41] but applies it specifically to AI's unique risk profile: the combination of speed, scale, and unpredictability.

**Practical Application** To apply this principle, organizations must shift from an "iterate in production" mentality to a "validate then scale" approach. This means defining explicit operational boundaries before deployment: What scenarios has this AI been trained and tested on? What inputs should trigger immediate human review? What statistical

anomalies indicate the AI is operating outside its validated domain? These boundaries must be encoded as hard constraints, not suggestions. For example, if a fraud detection AI starts flagging more than X% of transactions in a given hour—significantly deviating from baseline—it should automatically throttle down and alert humans, rather than continuing to execute. The most common pitfall is overconfidence from initial testing. Just because an AI performed well in a controlled environment doesn't mean it will perform well in production, where edge cases and novel scenarios are guaranteed. Implementation requires building extensive simulation environments—digital twins of real-world operations—where AI systems can encounter millions of edge cases before ever touching real decisions. The goal is to discover failure modes at simulation speed, not production speed.

# ***People Own Objectives***

Every objective needs a human owner to ensure people remain accountable for outcomes. When AI causes harm, the human owner is accountable, not the algorithm. _Name the Owner._

**The Hidden Problem** The arrival of sophisticated AI has introduced a dangerous new entity into our accountability frameworks: the "black box." When an AI system produces a harmful, biased, or simply incorrect outcome, the first human instinct is often to point a finger at the algorithm itself. This creates a tempting and perilous illusion of diffused responsibility, where a non-sentient tool is treated as a culpable actor. This phenomenon, often termed "automation bias," describes our tendency to over-trust automated systems and abdicate our own judgment, a cognitive shortcut that becomes exponentially more dangerous as those systems grow in complexity.[^17] The hidden problem is not that AI will make mistakes; it is that organizations will create structures that allow those mistakes to become nobody's fault. This accountability vacuum allows for systemic risks to grow unchecked, as no single individual feels the ownership necessary to question, audit, or challenge the machine's output. The system becomes an orphan, and its failures become organizational acts of God.

**The Cost of Ignoring It** When no one person owns the objective, the AI system becomes a scapegoat for organizational failure. A biased lending algorithm that denies loans to qualified applicants is not a "coding error"; it is a failure of the human who owns the objective of fair lending. A supply chain model that collapses under unexpected disruption is not a "model drift" problem; it is a failure of the human who owns the objective of resilient logistics. Ignoring this principle leads to a corrosive culture of finger-pointing and risk aversion. Instead of building robust systems, teams build elaborate justifications. Meetings are held not to solve problems, but to document that one is not to blame for them. This behavior is a direct echo of what Fred Brooks identified in *The Mythical Man-Month*: the diffusion of responsibility across a team leads to a state where no single person can be held accountable, and the project itself suffers.[^18] In the age of AI, this cost is magnified. A single algorithmic error can

propagate across millions of decisions in an instant, making the absence of clear ownership an existential threat.

**The Contrarian Insight** The conventional approach to AI safety focuses heavily on technical solutions: explainable AI (XAI), model auditing, and bias detection tools. While essential, these are insufficient. The contrarian insight of this principle is that the root of AI safety is not technical, but organizational. Accountability is a human construct, and it cannot be delegated to silicon. The most powerful tool for ensuring a safe and effective AI system is a single human name assigned to its objective. This is not about blaming an individual; it is about empowering one. The designated owner is not the person who writes the code or configures the model. They are the person who is accountable for the *business outcome* the AI is intended to achieve. This shifts the focus from "Is the model accurate?" to "Is the objective being met safely and effectively?" This single point of ownership creates the moral and professional gravity necessary to force hard questions and demand true understanding, transforming accountability from a distributed, abstract concept into a concrete, human responsibility.

**Solution Framework** Implementing this principle requires a simple but radical act: for every AI system deployed, an executive "owner" must be publicly named. This individual is accountable for the system's outcomes, good or bad. This framework forces a cascade of positive behaviors. The owner, now personally on the hook, is intensely motivated to understand the system's limitations. They are compelled to ask "how could this go wrong?" and to demand safeguards that go beyond technical minimums. They become the chief skeptic, challenging the enthusiastic claims of the development team. This approach mirrors the principles of Total Quality Management (TQM), where responsibility for quality is not relegated to a downstream inspection department but is owned by those who manage the process itself.[^19] Naming the owner makes accountability a proactive, design-time consideration, not a reactive, post-mortem exercise. It forces the organization to move beyond plausible deniability and embrace genuine responsibility.

**Philosophical Foundation** This principle is deeply rooted in the philosophy of human-computer interaction and organizational psychology. Decades of research have shown that technology is not a neutral actor; it is a medium that reflects and amplifies the values of its creators.[^20] This connects to the foundational concept of the "sociotechnical system," first identified by researchers at the Tavistock Institute, which posits that organizational outcomes are the product of inseparable interactions between social structures and technical tools.[^21] A system without a human owner is a technical object detached from the social contract of responsibility. Furthermore, this principle aligns with the core tenets of Servant Leadership, where the leader's primary role is to ensure the team is functioning effectively and ethically toward a clear goal. [^22] The AI owner is a servant leader for the human-AI system, responsible for its well-being and its impact on others. Finally, the act of naming an owner is a direct

application of the Agile Manifesto's emphasis on "individuals and interactions over processes and tools."[^23] It correctly identifies that the ultimate arbiter of a system's success is not the tool itself, but the human judgment guiding it.

Implementation Reality In practice, naming a single owner will face immediate and fierce resistance. It cuts against the grain of corporate cultures that favor committee-based decisions and shared responsibility precisely because those structures diffuse risk and protect individuals from blame. The first pushback will be an attempt to name a committee or a department as the owner; this must be rejected. A committee has no single throat to choke. The second point of resistance will come from potential owners who fear becoming a scapegoat for complex system failures. This fear is valid and must be addressed by framing the role as one of empowerment, not blame. The owner must be given the ultimate authority to question, halt, or modify the AI's operation to meet their objective. They are not just responsible; they are in charge. Overcoming this resistance requires unambiguous executive sponsorship and a clear, organization-wide understanding that the goal is not to punish failure, but to create the conditions for success.

# ***Deception Destroys Trust***

AI that pretends to be human eliminates informed consent and creates false relationships. People cannot collaborate effectively with what they don't recognize as artificial. *Make AI obvious, not hidden.*

The Operational Dysfunction There is a growing trend to design AI systems, particularly chatbots and voice assistants, to be as "human-like" as possible. They are given names, scripted with conversational tics, and designed to project personality and even emotion. The intention may be to create a more "natural" or "friendly" user experience. However, this creates a profound operational chaos: it introduces deception into the very foundation of the user interaction. When a person does not know they are talking to a machine, they cannot give informed consent to the interaction. They may disclose more information than they otherwise would, or form a one-sided emotional connection —a parasocial relationship—with a system that cannot reciprocate. This deception is not a harmless gimmick; it is a violation of the user's autonomy and a corrosive agent that makes genuine trust impossible.

Why Standard Approaches Fail The standard approach is driven by a flawed metric: "naturalness." Design teams strive to pass the Turing Test, not because it creates a better outcome for the user, but because it is a technical challenge. This focus on biomimicry is fundamentally misguided. The goal of a collaborative system should be clarity and effectiveness, not impersonation. A hammer is a good tool not because it looks like a human fist, but because its form perfectly communicates its function. Standard approaches fail because they conflate "easy to use" with "human-like." As Don Norman argues in *The Design of Everyday Things*, the best designs are those

that provide clear "signifiers" that communicate how an object can be used.[^28] An AI that pretends to be human provides a false signifier. It suggests a capacity for empathy, understanding, and reciprocity that simply does not exist, setting the user up for disappointment and a feeling of being duped.

**The Principle's Core Logic** The core logic of this principle is that collaboration requires trust, and trust requires honesty. You cannot have a healthy, functional collaboration with a partner who is lying to you about their fundamental identity. Making AI obvious is not about creating cold, robotic interfaces. It is about establishing clear boundaries and managing expectations. When an AI clearly identifies itself as such—"I'm the company's automated scheduling assistant"—it frames the interaction correctly. The user understands they are interacting with a tool, not a person. They will adjust their language, their expectations, and their level of disclosure accordingly. This transparency creates a foundation of trust. The user trusts that the system will be good at its stated purpose (scheduling a meeting) and will not pretend to be good at something it is not (offering emotional support). This honesty, as counterintuitive as it may seem to some marketers, is the only sustainable path to long-term user adoption and engagement. Deception is a short-term trick; transparency is a long-term strategy.

**Research Validation** This principle is strongly supported by research in HCI and technology ethics. Studies on anthropomorphism show that while users may initially find human-like agents appealing, this effect quickly sours when the agent fails to meet the heightened social expectations it created, leading to a greater sense of betrayal and frustration.[^33] This is related to the "Uncanny Valley" hypothesis, which suggests that near-perfect human replicas can elicit feelings of unease or revulsion.[^34] Foundational ethical frameworks prioritize informed consent as a cornerstone of any interaction, be it medical, financial, or digital.[^35] An AI that hides its nature denies the user the ability to consent to the terms of the interaction. Further research in social psychology confirms that perceived dishonesty is one of the fastest and most potent destroyers of trust in any relationship.[^36] Making AI obvious respects the user's cognitive autonomy, establishing a partnership based on clarity rather than an illusion.

**Practical Application** Practically, this means establishing clear organizational design standards that forbid deceptive practices. Chatbots should have names that are clearly non-human (e.g., "SchedulingBot 3000") or must explicitly state their nature at the beginning of an interaction ("Hi, I'm a virtual assistant..."). AI-generated text, images, or voice should be watermarked or otherwise labeled as synthetic. The goal is to create an information-rich environment where the user is never in doubt about the nature of the entity they are dealing with. The most common pitfall is the belief that this transparency will make the product seem less advanced or "magical." This is a failure of confidence in the product's actual utility. If an AI tool is genuinely useful, it does not need to pretend to be human. Its value will be self-evident. The focus should shift from "human-like" to "hyper-competent." Make the AI so good at its job that users are delighted to use it, knowing full well it is a machine.

### ***Individuals First***

AI industrializes manipulation by personalizing it at scale. Build tools that people control, not tools that control people. _Prioritize individual agency above efficiency, profit, or convenience._

**The Hidden Problem** The dominant logic of technology development is optimization. We are conditioned to seek efficiency, to minimize friction, and to maximize engagement or output. When applied to AI, this logic carries a hidden and dangerous payload. What makes modern AI uniquely powerful—its ability to personalize at scale—also makes it uniquely dangerous. An AI designed to maximize profit for an insurance company may learn that denying claims to the most vulnerable customers is the most effective strategy. An AI designed to maximize employee productivity may conclude that constant surveillance and algorithmically generated pressure are the best tools. The hidden problem is that AI doesn't just automate manipulation; it industrializes it by tailoring persuasion to individual psychological vulnerabilities. What once required expensive mass campaigns—focus groups, demographic targeting, broad messaging—now happens at the individual level, operating faster than people can recognize or meaningfully consent to. By prioritizing abstract metrics over human well-being, we risk building systems that are exquisitely effective at achieving the wrong goals, but do so with precision that makes resistance nearly impossible.

**The Cost of Ignoring It** When human agency is subordinated to system goals, the result is a landscape of digital puppetry at industrial scale. We see this in dark patterns, where user interfaces are intentionally designed to trick people into making choices they would not otherwise make, a direct assault on autonomy.[^24] We see it in social media recommendation engines that learn to serve increasingly extreme content because engagement metrics reward emotional manipulation.[^25] We see it in gig economy platforms where algorithmic managers dictate every action, leaving workers with a profound sense of powerlessness and alienation.[^26] The cost is not just a loss of trust; it is a measurable decline in mental and emotional well-being. This aligns with research on psychological safety, which demonstrates that environments where individuals feel controlled or threatened suffer from a catastrophic drop in innovation, engagement, and honest feedback.[^27] An organization that ignores this principle will inevitably create AI systems that, while perhaps profitable in the short term, generate a long-term trust deficit with its customers and a culture of fear among its employees. They will build systems that are technically efficient but morally and operationally bankrupt.

**The Contrarian Insight** The conventional wisdom is that a trade-off exists between human-centric design and business objectives like profit or efficiency. The contrarian insight is that this is a false dichotomy. In the long run, systems that undermine human autonomy are brittle and unsustainable. They generate resentment, workarounds, and eventual abandonment. The most robust and profitable systems are those that empower users, not manipulate them. This principle asserts that human well-being is not a constraint on design, but the ultimate objective of it. A system that makes a user

feel respected, capable, and in control is a system that will earn their loyalty and engagement. This shifts the design question from "How can we get the user to do what we want?" to "How can we help the user achieve their goal with dignity and clarity?" This is the core philosophy of human-centered design, which argues that the most successful products emerge from a deep and empathetic understanding of the user's needs and context, not just the organization's.[^28] The unique danger of AI is that it can execute manipulation at a scale and speed that makes individual resistance futile, transforming what were once isolated dark patterns into an ambient environment of algorithmic coercion.

**Solution Framework** To put individuals first, organizations must embed the concept of "dignity" as a primary design constraint. This means that during the design and review process, teams must be able to answer a critical question: "Does this system treat the human interacting with it as an end, or as a means to an end?" This framework requires the creation of "dignity metrics" that sit alongside performance metrics. For example, does the system offer clear and easy ways to opt-out or appeal a decision? Does it provide transparency about how it works, or does it operate as an unexplainable black box? Does it create pressure or provide support? Does it adapt to individual vulnerabilities to manipulate behavior, or does it provide clear choices? This approach involves explicitly mapping the emotional and psychological journey of the user, not just their transactional path. It requires adopting techniques like "value stream mapping" not just for process efficiency, but for human experience, identifying and eliminating steps that cause frustration, confusion, or a sense of powerlessness.[^29]

**Philosophical Foundation** This principle stands on the shoulders of centuries of philosophical thought regarding human autonomy, most famously articulated by Immanuel Kant's categorical imperative to always treat humanity as an end, and never merely as a means.[^30] In the age of AI, this is no longer an abstract philosophical exercise; it is a concrete design principle for preventing instrumentalization, where a person is reduced to a mere tool for the system's goal. Technology that uses deception or manipulation to achieve its ends explicitly violates this imperative. This principle also draws heavily from the field of positive computing, which focuses on designing technology to support human well-being and flourishing, arguing that our tools should not just be functional but should contribute positively to our psychological state.[^31] It asserts that good technology should enhance human capability and experience, not diminish it for the sake of an optimized metric. By placing the individual first, we commit to building AI that serves human values, rather than asking humans to serve the values of the machine. The AI-specific danger is scale: what was once a local manipulation by a salesperson becomes an ambient, inescapable environment of algorithmic influence operating on billions of people simultaneously.

**Implementation Reality** The primary resistance to this principle will be the perceived cost and the pressure of metrics. Building in transparency, offering recourse, and prioritizing user control can seem more expensive and less "efficient" than creating a

closed, optimized system designed to maximize a single KPI like conversion or engagement. The argument will be made that "users don't care" or "it will hurt our numbers." This resistance is best overcome with evidence, not argument. Start with small, contained projects that demonstrate the long-term value of a trust-based relationship. Show how empowered users become more loyal customers and how respected employees are more innovative. This requires a willingness to measure success in years, not quarters, and to believe that the most sustainable competitive advantage is the trust of the people you serve. It is a commitment to building systems that we would, ourselves, be willing to be subject to. The escalating concern with AI is that the optimization loop operates faster than human awareness, making it critical to embed dignity constraints at the architectural level rather than hoping humans can resist individually.

# ***Build From User Experience***

Design insight comes from living with the daily friction that analysis misses. People who navigate these daily realities understand what breaks and why. _People wrestling with system failures are the ones qualified to design system futures._

The Hidden Problem There is a fundamental disconnect in how most systems are designed. The people who conceive of the system (executives and architects) are institutionally insulated from the consequences of its daily use. They see the system through flowcharts, spreadsheets, and high-level dashboards—clean, abstract representations that bear little resemblance to the messy reality of its operation. The people who actually use the system—the customer service agent, the factory floor worker, the end customer—experience it as a series of frustrating workarounds, confusing interfaces, and illogical dead ends. The hidden problem is that organizations systematically devalue the most crucial source of design insight: the lived experience of failure. This experiential knowledge, rich with context and nuance, is often dismissed as "anecdotal" in favor of aggregated quantitative data that masks the very friction points that signal deep design flaws. When AI is designed from the boardroom rather than the front lines, it automates the wrong things—codifying the broken process rather than fixing the underlying dysfunction.

The Cost of Ignoring It When you design from the top down, you build elegant solutions to the wrong problems. This is the path to failed products, wasted investment, and user resentment. An organization might spend millions on a new AI-powered enterprise resource planning (ERP) system designed to "streamline" operations, only to find that employees have developed a complex shadow system of spreadsheets because the official tool is unusable for their actual tasks. This is a classic example of what Peter Senge calls "compensating feedback," where a well-intentioned intervention is defeated by the system's own response to it.[^48] The cost is not just the wasted money; it is the erosion of trust between leadership and the front lines. Employees become cynical, believing that their real-world expertise is ignored. Customers

become frustrated, feeling that the company does not understand or care about their experience. This creates a cycle where bad systems are endlessly patched but never fundamentally fixed. With AI, this problem accelerates: a poorly designed system doesn't just create friction—it scales that friction to millions of interactions before anyone realizes the core assumption was wrong.

**The Contrarian Insight** The contrarian insight, borrowed from the world of human-centered design, is that the person experiencing the problem is more of an expert than the person analyzing it from a distance. While data can tell you *what* is happening, only lived experience can tell you *why*. This principle argues that the most valuable design activity is not a brainstorming session in a conference room, but an ethnographic immersion into the user's world. It is about understanding the "job to be done" from the user's perspective, a concept powerfully articulated by Clayton Christensen.[^49] He argued that people don't buy products; they "hire" them to do a job. To design a better product, you must deeply understand the job, not the customer demographics. The people on the front lines, wrestling with the system's failures, are the ones who understand the "job" most intimately. They are the true subject matter experts. For AI systems, this is especially critical: if you don't understand the real-world context in which decisions are made, you'll train AI on sanitized data that doesn't reflect actual conditions, building systems that fail when they encounter reality.

**Solution Framework** To build from user experience, organizations must invert the traditional design hierarchy. Instead of architects handing down blueprints, designers must become embedded observers and co-creators with users. The framework for this is grounded in the principles of Design Thinking, which begins with empathy.[^50] This means developers, product managers, and even executives must spend meaningful time "on the floor," directly observing and even performing the work their systems are meant to support. This practice, known in Lean manufacturing as *genchi genbutsu* ("go and see"), is foundational to understanding the reality of a problem.[^51] The insights from this direct observation then become the raw material for ideation and prototyping. The goal is not to ask users what they want—they often cannot articulate it —but to understand their context so deeply that you can design a solution that feels like it was built just for them. For AI specifically, this means observing not just the happy path, but the exceptions, edge cases, and workarounds that reveal where current systems fail.

**Philosophical Foundation** This principle is rooted in the philosophy of phenomenology, which posits that subjective experience is a primary source of knowledge that cannot be fully captured by objective measurement.[^52] It is also a direct application of systems thinking, which teaches that a system's behavior is an emergent property of the interactions between its parts—a property that cannot be understood by analyzing the parts in isolation.[^53] A user's frustration is an emergent property of a poorly designed system, and it can only be diagnosed by observing the user interacting with the whole. This principle also embodies the core ideas of the Agile Manifesto,

particularly the value of "customer collaboration over contract negotiation."[^54] It extends this idea from the paying customer to all users of a system, internal and external, treating them as essential partners in the creative process. It is a commitment to intellectual humility—an acknowledgment that true expertise is often found farthest from the executive suite and closest to the work itself.

Implementation Reality The greatest barrier to implementing this principle is organizational ego and the illusion of efficiency. It requires leaders to admit they do not have all the answers and that the most valuable insights may come from the lowest-paid employees. It also requires an investment of time that feels "unproductive" to managers focused on immediate output. Sending a team to do ethnographic research feels slower than jumping straight into development. This resistance must be countered by demonstrating the high cost of *not* doing it: the failed projects, the wasted engineering cycles, and the abandoned products. The most effective tactic is to deliver a rapid prototype that solves a real, observed point of friction. When an executive sees a simple solution that perfectly addresses a problem they didn't even know existed, the value of the approach becomes undeniable, proving that the shortest path to a successful product is not a straight line from idea to launch, but an iterative loop that begins and ends with the human experience.

## ***Discovery Before Disruption***

Systems hide logic until something breaks. The redundancy that looks pointless prevents failures you've never seen. The manual step that feels inefficient satisfies requirements nobody documented. Deletion scales faster than comprehension. _Remove only what you understand; build to discover the rest._

The Operational Dysfunction Driven by an eagerness to innovate, technical teams often approach a legacy system with the primary goal of replacing it. They see an old, clunky process and their first instinct is to tear it down and build something new, elegant, and modern from scratch. This "rip and replace" mentality is incredibly risky. It is born of arrogance—the belief that the original designers were unsophisticated and that the current users are simply tolerating a bad system. What this approach fails to recognize is that long-standing systems have often evolved complex, invisible mechanisms to handle undocumented exceptions and edge cases. The seemingly "pointless" manual review step might be there to catch a specific type of fraud that the formal rules don't account for. The redundant data entry might be the only mechanism that reconciles discrepancies between two legacy databases. When AI is deployed to "optimize" these processes, it moves at machine speed—deleting what looks inefficient before anyone realizes those inefficiencies were actually critical safeguards. Deletion scales faster than comprehension.

Why Standard Approaches Fail Standard approaches to system analysis focus on the documented, official workflow—the "happy path." They model the process as it is

*supposed* to work. This fails because the most important logic in a human system often exists in the undocumented workarounds, the informal networks, and the unwritten rules that people use to make the official process actually function. This is the essence of Chesterton's Fence principle: do not tear down a fence until you understand why it was put up in the first place.[^43] Standard analysis looks at the fence and sees an obstacle to a clean, open field. It fails to investigate whether that fence is protecting the field from something it cannot see. By ignoring the "as-is" reality in favor of a theoretical "to-be" design, these approaches set the stage for building a beautiful new AI system that fails catastrophically as soon as it encounters the messy realities of the real world. With AI, this failure happens faster and at greater scale than with traditional systems because AI doesn't pause to ask questions—it executes its interpretation immediately and continuously.

**The Principle's Core Logic** The core logic of this principle is one of institutional humility. It insists that before you earn the right to disrupt a system, you must first earn a deep and empathetic understanding of it. This means treating the existing process not as a problem to be solved, but as a source of invaluable data. The goal of the initial discovery phase is not to design the new system, but to become an archaeologist of the old one. Why do people do what they do? What are the hidden pressures and incentives? What crises has this system survived, and what adaptations did it make? This is about seeking first to understand, then to be understood. Only after you have mapped the hidden logic, identified the unwritten rules, and understood the "why" behind every seemingly illogical step can you begin to design a replacement that is not just more efficient, but more resilient. With AI specifically, this means understanding not just the current process, but the full range of edge cases and exceptions that the AI will need to handle—because unlike humans, AI won't improvise when it encounters something unexpected.

**Research Validation** This principle draws heavily from the fields of ethnography and systems thinking. Ethnographic research methods, often used in anthropology, involve immersing oneself in a culture to understand its implicit rules and behaviors—a perfect analogy for understanding a complex legacy system where the "official" rules rarely tell the whole story.[^56] It is also a direct application of Peter Senge's "Laws of the Fifth Discipline." One of Senge's key laws is that "behavior grows better before it grows worse," warning that a quick, superficial fix often creates short-term improvement but long-term disaster because it disrupts hidden systemic balances that were never understood.[^57] This principle operationalizes the core Lean principle of *genchi genbutsu* ("go and see"), which commands that one must go to the source of the work to truly understand it, rather than analyzing it from a distance.[^58] In the context of AI, this discovery phase is critical because once an AI system is deployed, it will execute its learned behavior at scale—if that behavior is based on misunderstanding the system's true requirements, the damage compounds rapidly.

Practical Application Practically, this means the first phase of any AI transformation project should be purely observational and anthropological. The team's job is to build a detailed "map" of the existing territory, including all the informal workarounds, shadow IT, and hidden social structures. A powerful tool for this is value stream mapping, but one that is focused on information flow and exception handling, not just the official process steps.[^59] The team should actively seek out the "gurus"—the long-time employees who know all the system's secrets—and treat them as revered sources of wisdom, not as relics of an old way of working. For AI implementation, this means creating exhaustive documentation of edge cases and exceptions before ever training a model—because the AI will need explicit rules for handling scenarios that humans currently manage through intuition and institutional knowledge. A common pitfall is for this discovery phase to be seen as a delay to the "real work" of building. Leadership must champion this phase as the most critical part of risk mitigation. The map you create during discovery is the single most important tool for ensuring your new AI system doesn't execute thousands of flawed decisions before you realize you misunderstood a critical requirement.

## ***Ambiguity Is Wisdom***

Experts navigate gray areas; beginners demand binary answers. AI produces probabilities that demand judgment, not facts that replace it. When systems force ambiguity into yes/no decisions, they destroy the space where expertise operates. _Reveal the probabilities._

The Operational Dysfunction There is a deep-seated human bias for certainty. We want clear, definitive answers, especially from our technology. When asked "Will this customer churn?" or "Is this transaction fraudulent?", we want a simple "yes" or "no." This creates immense pressure on developers to design AI systems that provide the illusion of certainty, even when the underlying reality is probabilistic. An operational breakdown occurs when a nuanced, probabilistic output (e.g., "85% confidence this is fraud") is converted into a simple binary answer—a process that destroys crucial information. A system that simply flags a transaction as "fraud" gives the human reviewer no context about the strength of that assessment. A system that says "85% confidence" gives the reviewer a powerful piece of information that can guide the depth and urgency of their investigation. The dysfunction emerges because binary decisions are organizationally convenient: they're easy to report, easy to act on, and easy to defend. But this convenience comes at the cost of wisdom. By hiding the "maybe," we build systems that are simultaneously dumber and more dangerously arrogant—they make sweeping pronouncements without revealing their uncertainty.

Why Standard Approaches Fail Standard approaches to system design optimize for simplicity and speed of decision-making. A manager wants a clean dashboard with red, yellow, and green lights, not a complex scatter plot of probabilities. This desire for what Daniel Kahneman calls "cognitive ease" leads us to build systems that cater to our biases rather than challenge them.[^42] We design systems that give us the

simple answers we crave, even if they are misleading. This fails because the real world is messy, complex, and inherently probabilistic. A medical diagnosis AI that gives a binary answer is useless; a doctor needs to understand the confidence level, the potential differential diagnoses, and the factors driving the model's conclusion. By oversimplifying the output for the sake of a clean user interface, standard approaches strip away the very nuance that makes an AI's prediction useful to an expert. They treat the human user as a passive recipient of an answer, rather than as an active partner in a process of sensemaking. The problem isn't uncertainty, it's active ambiguity—the recognition that real-world situations don't fit neatly into categories, and that the ability to work in gray areas is what differentiates experts from beginners.

**The Principle's Core Logic** The core logic of this principle is that uncertainty is not noise; it is data. A probability score, a confidence interval, or a range of possible outcomes contains valuable information about the reliability of a prediction. Presenting this ambiguity to the user is an act of intellectual honesty and a prerequisite for effective human-AI collaboration. It allows the human expert to apply their own contextual knowledge to the AI's probabilistic output. If the AI is 99% confident, the human can proceed quickly. If the AI is 60% confident, the human knows to slow down, gather more information, and apply a higher degree of skepticism. This approach transforms the AI from an oracle that provides answers into a tool that provides evidence. It empowers the user, respecting their expertise and giving them the information they need to make a better final judgment. When systems force ambiguity into binary decisions, they destroy the space where expertise operates—the gray areas where judgment, context, and experience matter most. This is about designing for wisdom, not just for answers.

**Research Validation** This principle is deeply rooted in the fields of statistics and decision theory. The entire framework of Bayesian inference, for example, is built around the idea of updating our beliefs in light of new, uncertain evidence, rather than seeking absolute truth.[^44] Forcing a probabilistic world into a deterministic model is a statistical sin that destroys information. This principle is also supported by extensive research in sensemaking, a concept pioneered by organizational theorist Karl Weick. Weick argued that people and organizations are constantly trying to make sense of ambiguous environments, and that the process of "sensemaking" is about creating plausible accounts of what is happening, not about finding a single, objective truth. [^45] An AI that presents uncertainty is providing the raw materials for sensemaking. It gives the user clues, probabilities, and ranges that they can weave into a more resilient and nuanced understanding of the situation, rather than a single, brittle "fact" that may be dangerously wrong. Research on expert decision-making consistently shows that experts don't seek binary answers; they seek information that helps them understand the situation's complexity.[^46]

**Practical Application** In practice, this means rejecting binary outputs for any non-trivial prediction. Instead of a "yes/no" flag, the system should display a confidence score.

Instead of a single predicted sales number, it should show a probable range. The user interface must be designed to visualize uncertainty effectively, using techniques like error bars, probability distributions, or even verbal framing like "likely" vs. "very likely." A common pitfall is information overload. Simply dumping a wall of statistical data on a user is not helpful. The key, as visualization expert Edward Tufte has long advocated, is to present complex information with clarity, precision, and efficiency, making it easy to understand at a glance.[^47] The goal is not to make the interface more complicated, but to make it more honest. It is a design challenge to represent nuance in a way that is intuitive and actionable, but it is a challenge that must be met if we are to build AI systems that truly make us smarter. The shift is from systems that pretend to know, to systems that reveal what they actually know—and what they don't.

## ***Reveal the Invisible***

Gaps in understanding hide inside abstraction until forced into concrete form. The most valuable representation is whatever hurts most to produce; whether diagram, specification, or working prototype. Easy articulation reveals nothing; difficulty exposes confusion. _Choose representations that force confrontation with what you don't know._

> The Operational Dysfunction Organizations are drowning in text. We communicate about complex systems through dense documents, lengthy email chains, and bullet-pointed slide decks. This reliance on prose to describe dynamic, multi-dimensional processes is a massive operational obstacle. Written language is linear and sequential, while most organizational processes are parallel, interconnected, and cyclical. More insidiously, text allows teams to maintain the comfortable illusion of shared understanding when none exists. Everyone nods along in the meeting, confident they understand the plan, only to discover months later that each person had a completely different mental model. The operational dysfunction is that abstraction hides confusion. You can write thousands of words about how an AI system should work without ever confronting the specific logical contradictions or edge cases that will cause it to fail. The words create a fog of plausible ambiguity that allows critical misunderstandings to persist unchallenged.

> Why Standard Approaches Fail Standard approaches to documentation and communication prioritize comprehensive detail over shared understanding. The goal becomes producing a "complete" requirements document or project charter, often running hundreds of pages. This fails for a simple human reason: nobody reads it. And even if they did, the text-based format makes it impossible to see the system as a whole. It is like trying to understand a city by reading a list of all its street names. You have no sense of the layout, the neighborhoods, the traffic flows, or the relationships between the parts. You can hide a wealth of ignorance inside a couple of very pretty diagrams—and you can hide even more inside pages of prose.[^60] The failure is that these approaches optimize for the appearance of rigor without forcing the

confrontation with complexity that reveals gaps in thinking. Standard documentation lets you describe complexity without demonstrating you understand it.

**The Principle's Core Logic** The core logic of this principle is that shared understanding requires forced confrontation with gaps in knowledge. The most valuable representation is not necessarily a visual one—it's whatever hurts most to produce. [^60] If you're claiming to understand a process, can you diagram it? If you're proposing an algorithm, can you write pseudocode for it? If you're defining a user experience, can you build a clickable prototype? The medium matters less than the forcing function: any representation that makes your abstract claims concrete will immediately expose where your thinking is fuzzy. A working prototype reveals assumptions you didn't know you were making. A detailed system diagram surfaces conflicts between components you thought were compatible. A journey map exposes moments of user confusion you never anticipated. The goal is not to create a beautiful artifact, but to force yourself and your team to confront what you don't actually know. Easy articulation reveals nothing; difficulty exposes confusion.

**Research Validation** This principle is validated by decades of research in cognitive psychology and information visualization. Cognitive load theory suggests that our working memory is extremely limited, and well-chosen external representations can offload complex information, making it easier to process and understand than purely verbal or textual explanations.[^61] The work of Edward Tufte provides a foundational framework for how to display quantitative and qualitative information with clarity and integrity, arguing that good design reveals the truth of the data, not just decorates a page.[^62] This principle is also a core component of Design Thinking, which relies heavily on tangible artifacts like storyboards, sketches, and prototypes to explore and communicate ideas because they are faster to create, easier to critique, and more effective at eliciting honest feedback than abstract documents.[^63] Research on "boundary objects"—shared artifacts that different groups can use to coordinate and communicate—shows that concrete representations create alignment precisely because they force specificity.[^64] The mantra is "show, don't tell," because showing bypasses the ambiguity of language and creates a direct, shared context.

**Practical Application** Applying this principle means shifting the default mode of communication from writing to making. Before writing a dense requirements document, create a prototype. Before debating system architecture in prose, draw it. Before claiming to understand a user's journey, map it. The specific medium matters less than the commitment: choose whatever representation will most painfully expose gaps in your current understanding. For AI systems specifically, this means creating concrete examples of edge cases and decision trees before claiming the AI can handle them. A common pitfall is "analysis paralysis," where teams spend too much time perfecting the representation instead of using it to discover what they don't know. The artifact is a disposable tool for thinking, not a deliverable. The moment the representation stops generating productive conversation or revealing new insights, its primary job is done.

For AI, forcing concreteness early—through simulation, prototyping, or scenario modeling—reveals failure modes while they're still cheap to fix, rather than discovering them in production when they've already been executed thousands of times.

## ***Iterate Towards What Works***

Requirements emerge through building, not planning meetings. Inherited practices carry outdated logic that meetings can't expose. Iteration without feedback is repetition; only rapid cycles of making, testing, and failing reveal what actually works. *Build to discover; test to validate; repeat.*

**The Operational Dysfunction** Organizations have a deep-seated belief that they can plan their way to success. This manifests in the creation of massive, detailed project plans and requirements documents, often developed over months of meetings, before a single line of code is written. The underlying assumption is that it is possible to fully understand and specify a complex system in advance. The root of the operational issue is this "waterfall" mindset, which frames building as the execution of a predefined plan. In reality, for any novel or complex problem, the plan is a hypothesis, and the act of building is the experiment that tests it. By front-loading all the "thinking" work, organizations create rigid plans based on flawed assumptions and ensure they will learn about their mistakes only when it is most expensive to fix them. With AI, this dysfunction is magnified: you cannot predict how an AI will behave across all scenarios through planning alone. You discover its failure modes by testing it in conditions that approximate reality.

**Why Standard Approaches Fail** Standard waterfall-style project management fails because it is fundamentally at odds with the nature of discovery and learning. It assumes a linear path from problem to solution in a world that is non-linear and unpredictable. As Fred Brooks noted in his seminal essay "No Silver Bullet," the hardest part of building software is not the coding, but the conceptual design—the formulation of the complex, abstract requirements.[^65] This conceptual work cannot be perfected on paper. Real understanding only emerges when you try to translate the abstract idea into a concrete artifact. It is only when you put a prototype in front of a real user, or test an AI on real data, that you discover the fatal flaw in your core assumption. Standard approaches that delay this moment of truth for as long as possible are optimizing for perceived predictability at the cost of actual success. The illusion is that comprehensive planning reduces risk; the reality is that it just delays the discovery of risk until the stakes are highest.

**The Principle's Core Logic** The core logic of this principle is that for complex systems—especially AI systems—building *is* thinking. The fastest way to learn is to create something tangible, expose it to reality, and see how reality responds. This is the foundational principle of the Agile Manifesto, which values "working software over comprehensive documentation" and "responding to change over following a plan." [^66] But iteration alone is insufficient. As the principle states: iteration without

feedback is mere repetition. The goal is to create a rapid feedback loop where each cycle of making and testing produces validated learning that informs the next cycle. This means deliberately designing small experiments that can fail safely and teach you something concrete. For AI specifically, this means building simulation environments where you can test the AI against thousands of scenarios before deploying it—discovering failure modes at simulation speed rather than production speed. The principle recognizes that inherited practices carry outdated logic that meetings can't expose; only the act of building and testing forces confrontation with reality.

Research Validation This principle is validated by decades of experience in product development and is supported by theories of organizational learning. Peter Senge's work on learning organizations emphasizes that genuine learning happens through action and reflection, not through abstract planning.[^67] The process of iterative prototyping is a powerful engine for this kind of learning, as it makes abstract ideas concrete and creates opportunities for rapid feedback. This principle also reflects the scientific method itself: form a hypothesis (the requirement), conduct an experiment (build a prototype), analyze the results (get user feedback), and refine the hypothesis. It treats product development not as a manufacturing process with a fixed blueprint, but as a research process where the blueprint is the output, not the input. The Agile Manifesto provides the philosophical foundation, codifying these principles that were learned through hard experience in the software industry.[^68] For AI systems, the research on machine learning validation emphasizes that model performance cannot be predicted from theory alone—it must be empirically tested against real-world data and edge cases.[^69]

Practical Application Applying this principle means radically shrinking the planning horizon. Instead of a six-month plan, create a plan for the next two weeks. The goal of that two-week "sprint" is not just to produce features, but to answer a critical question or test a key assumption. The output of the sprint is not just code; it is validated learning that informs the next sprint. This requires a cultural shift away from viewing changing requirements as a sign of failure, and towards seeing them as a sign of progress. A change means you have learned something you did not know before. The most common pitfall is "fake agile," where teams use agile terminology (sprints, stand-ups) but still operate with a waterfall mindset, treating iterations as mini-waterfalls with no real feedback loop or willingness to pivot based on what's learned. True iteration requires the intellectual humility to accept that your initial plan is probably wrong, and the organizational courage to change it based on evidence. For AI, this means building continuous validation into the development process—testing not just model accuracy, but behavior in edge cases and failure modes—and being willing to fundamentally rethink the approach when testing reveals the initial assumptions were flawed.

# ***Decompose Incrementally***

Legacy systems carry too much technical debt to replace and are too brittle to automate. AI systems should allow isolated components to decompose naturally. *Dismantle legacy complexity piece-by-piece.*

**The Operational Dysfunction** When faced with a complex, messy, and unreliable legacy system, the most seductive idea for any engineering team is the "Big Rewrite." The promise is to throw away the tangled mess of old code and convoluted processes and replace it with a clean, modern, perfectly architected solution built from scratch. This desire is understandable, but it is also incredibly dangerous. The operational dysfunction is twofold. First, the Big Rewrite ignores the hidden value embedded in the old system—the accumulated lessons from handling thousands of edge cases, the subtle logic that prevents failures you've never consciously observed. Second, and more critically, it creates a single, catastrophic point of failure: if the new system doesn't work, you've destroyed the old one and have nothing to fall back on. This all-or-nothing approach is the opposite of resilience. Systems built this way are brittle—they cannot adapt incrementally to changing assumptions or requirements. When conditions change (and they always do), the entire system becomes obsolete, forcing another costly, risky big rewrite.

**Why Standard Approaches Fail** Standard approaches often present a false choice between two extremes: either continue to apply small, incremental patches to a failing system, or embark on a massive, multi-year project to replace it entirely. The first option leads to a slow, painful death by a thousand cuts. The second is a high-stakes gamble that, according to industry data on large software projects, is more likely to fail than to succeed.[^75] The Big Rewrite approach fails because it violates the principle of iterative feedback. It is a return to the worst aspects of waterfall development, with a massive, singular point of failure at the final launch. Standard approaches fail because they don't design for obsolescence—they assume the system will remain relevant indefinitely, which it never does. They optimize for initial elegance rather than long-term adaptability.

**The Principle's Core Logic** The core logic, informed directly by Dave Thomas's feedback, is that systems should be designed from the outset for piece-by-piece evolution rather than wholesale replacement. This means building modularity and optionality into the architecture: clear boundaries between components, well-defined interfaces, and the ability to swap out one piece without disrupting the whole. The Japanese concept of *Monozukuri*—the art of making things—embodies this philosophy: respect for craft includes designing things that can be maintained, adapted, and incrementally improved over time, not just things that are initially impressive.[^76] The principle recognizes that all assumptions eventually break. Technologies evolve, markets shift, regulations change. A system that can only adapt through total replacement is doomed to obsolescence. A system designed for incremental obsolescence—where individual components can be retired and replaced

as they age or as better alternatives emerge—is antifragile: it gains from the stress of changing conditions rather than being destroyed by them.[^77]

Research Validation This principle is a practical synthesis of several established frameworks. It draws from the Lean Startup's emphasis on validated learning and iterative improvement,[^78] but applies it to system architecture rather than just product features. It aligns with John Kotter's research on leading organizational change, which found that successful, large-scale change efforts almost always begin by generating and publicizing "short-term wins"—unambiguous, visible improvements that provide credibility and momentum for larger initiatives.[^79] Building for incremental obsolescence is a strategy for creating those short-term wins: each component you successfully replace proves your capability and builds trust. It also embodies Chesterton's Fence principle: by replacing one piece at a time, you're forced to understand how each component works before you remove it, reducing the risk of catastrophic errors.[^43] The architectural principle of "modularity" in software engineering—designing systems as independent, loosely coupled components—is the technical implementation of this philosophy.[^81]

Practical Application In practice, this means designing system architecture with explicit replacement boundaries. Before building, ask: "If we needed to replace just this one component three years from now, could we do it without disrupting the rest of the system?" This requires investing in clear interfaces, comprehensive documentation of dependencies, and modular design even when monolithic approaches might be faster initially. When organizations propose a Big Rewrite, the first question from leadership should be: "Can you show us three significant component-level improvements you've successfully made to the current system?" If the team cannot demonstrate capability at the small scale, they haven't earned the right to attempt transformation at the large scale. A common pitfall is that engineers may see incremental replacement as unglamorous compared to the excitement of building something entirely new. It is the role of leadership to reframe this work as what it is: the professional, low-risk path to genuine transformation. For AI systems specifically, this means designing so that you can swap out models, data sources, or decision logic components independently— allowing the system to evolve as better AI techniques emerge without requiring a complete rebuild.

# ***Justify Resource Consumption***

AI makes it trivially easy to waste resources. What costs pennies to create can cost millions to run. The resources remain real: energy, water, compute, time. *Optimize the ratio of value per resource spent.*

The Operational Dysfunction Organizations are adopting AI with a brute-force mentality, equating more computational power with better outcomes. This leads to a new form of digital waste: using massive, general-purpose AI models for tasks that require only a

fraction of their capability. A team might deploy a state-of-the-art large language model to perform simple classification, a task a much smaller, specialized model could handle faster and cheaper. This isn't just inefficient; it's a systemic misalignment. Traditional systems made waste visible—you could see the paperwork piling up, the redundant meetings on calendars, the idle capacity on factory floors. AI makes waste invisible: compute cycles consumed in distant data centers, energy costs hidden in cloud bills, latency measured in milliseconds that aggregate into user frustration. The operational dysfunction is that AI compounds small inefficiencies at scale, transforming minor design decisions into massive hidden costs before anyone realizes the system is hemorrhaging resources. Not all complexity delivers value, but AI makes it easy to scale complexity that delivers none.

**Why Standard Approaches Fail** Standard approaches are captive to the "bigger is better" narrative of the AI industry, which is locked in a marketing-driven race toward ever-larger models. This thinking trickles down, leading implementation teams to believe that using the most powerful model is a sign of sophistication. This fails because it optimizes for the wrong variable: raw capability instead of task-appropriate efficiency. Everything must justify itself through value delivered per resource spent. Standard approaches fail to account for the significant and often hidden costs of over-engineering. As journalist Karen Hao has documented, the environmental toll of training and running massive models—in terms of energy and water consumption—is staggering and growing exponentially.[^70] Standard approaches mistake theoretical power for practical value, leading to the deployment of systems that are intelligent but operationally sluggish, economically wasteful, and environmentally unsustainable. They also fail to distinguish between complexity that creates competitive advantage (sophisticated fraud detection) and complexity that just creates work (five-approval purchase processes).

**The Principle's Core Logic** The core logic is to make value per resource spent the organizing principle for all system design. This reframes the question from "Can we do this with AI?" to "Should we do this with AI, and if so, what is the minimum effective intelligence required?" This requires attacking waste in all its forms. First, computational waste: choosing smaller, specialized models over large, general-purpose ones whenever possible. Second, organizational waste: the bureaucratic complexity that adds control without adding value. Third, temporal waste: delays and handoffs that create friction. The principle recognizes that AI can amplify both good complexity (that creates advantage) and bad complexity (that creates work). The goal is not simplicity for its own sake—some problems require sophisticated solutions. The goal is ensuring that every unit of complexity, every compute cycle, every approval gate can justify its existence through the value it delivers. If it can't, it should be eliminated. This is about building systems that are economically and operationally sustainable, not just technically impressive.

**Research Validation** This principle extends the foundational Lean manufacturing concept of eliminating *muda* (waste) into digital and AI domains.[^71] While Lean traditionally focused on physical inventory and wait times, this principle identifies computational overkill and unnecessary complexity as critical forms of modern waste. The argument is supported by research on the resource consumption of large-scale AI, which provides a tangible, economic, and ethical imperative for efficiency.[^72] Furthermore, the core idea aligns with the Theory of Constraints, which posits that any system's output is limited by a single bottleneck.[^73] In an AI-enabled workflow, an oversized, high-latency model can easily become the new constraint that throttles the entire process, no matter how optimized other steps are. This principle also connects to the work on sustainable computing, which argues that the environmental cost of computation must be factored into system design, not treated as an externality.[^74]

**Practical Application** Applying this principle means making "value per resource spent" a visible, primary metric for all system design decisions. Before deploying an AI component, the team must ask: "What is the simplest, smallest, most efficient approach that can achieve the required quality bar?" This might mean choosing a fine-tuned small model over a massive proprietary API for a specific task. It requires measuring not just the accuracy of a solution, but its latency, cost-per-use, energy consumption, and organizational friction as key performance indicators. For organizational complexity, it means conducting regular "complexity audits": mapping which processes and approvals actually deliver value versus which exist purely to satisfy internal bureaucracy. A common pitfall is "resume-driven development," where engineers choose complex, powerful solutions because they're more interesting, not because they're more appropriate. Leadership must counter this by celebrating and rewarding efficiency—making the total cost of a solution, including its computational and organizational burden, as important as whether it works. The goal is to scale only what earns its cost: the things that genuinely create competitive advantage or eliminate meaningful user pain.

# Copyright Notice

Requests For Support: We encourage, but do not require, the following to support the AI First Principles project:

1. **Notification:** Let us know how and where you have used the AI First Principles by emailing info@aifirstprinciples.org. 2. **Contribution:** If you have suggestions for changes or improvements, please share them via our GitHub repository at [https://github.com/aifirstprinciples/AI-First-Principles] (https://github.com/aifirstprinciples/AI-First-Principles). 3. **Financial Support:** Consider supporting the promotion and development of AI First Principles through our Patreon at [https://www.patreon.com/c/AIFirstPrinciples] (https://www.patreon.com/c/AIFirstPrinciples).

# References

[1] Ransbotham, Sam, et al. "The Cultural Benefits of Artificial Intelligence in the Enterprise." *MIT Sloan Management Review* and *Boston Consulting Group*, May 2022.

[2] Senge, Peter M. *The Fifth Discipline: The Art & Practice of The Learning Organization*. Doubleday/Currency, 2006.

[3] Norman, Don. *The Design of Everyday Things: Revised and Expanded Edition*. Basic Books, 2013.

[4] Wilson, Robb, and Josh Tyson. *Age of Invisible Machines: A Practical Guide to Creating a Better, More Productive, and More Human World of Work*. 2nd ed. Wiley, 2025.

[5] Goldratt, Eliyahu M., and Jeff Cox. *The Goal: A Process of Ongoing Improvement*. North River Press, 2014.

[6] Christensen, Clayton M., et al. "Know Your Customers' 'Jobs to Be Done'." *Harvard Business Review*, September 2016.

[7] Based on Dave Thomas feedback session and analysis, 2024. Thomas is co-author of the Agile Manifesto and emphasized that AI's unpredictability aligns more with human behavior than engineering predictability.

[8] National Institute of Standards and Technology. "Artificial Intelligence Risk Management Framework (AI RMF 1.0)." NIST, Jan. 2023, [https://doi.org/10.6028/NIST.AI.100-1] (https://doi.org/10.6028/NIST.AI.100-1).

[9] O'Neil, Cathy. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown, 2016.

[10] Dastin, Jeffrey. "Amazon scraps secret AI recruiting tool that showed bias against women." *Reuters*, October 10, 2018.

[11] Noble, Safiya Umoja. *Algorithms of Oppression: How Search Engines Reinforce Racism*. NYU Press, 2018.

[12] Gebru, Timnit, et al. "Datasheets for Datasets." *Communications of the ACM*, vol. 64, no. 12, 2021, pp. 86-92.

[13] Raji, Inioluwa Deborah, et al. "Closing the AI Accountability Gap: Defining an End-to-End Framework for Internal Algorithmic Auditing." *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020, pp. 33-44.

[14] Haraway, Donna. "Situated Knowledges: The Science Question in Feminism and the Privilege of Partial Perspective." *Feminist Studies*, vol. 14, no. 3, 1988, pp. 575-599.

[15] Diakopoulos, Nicholas. "Algorithmic Accountability: Journalistic investigation of computational power structures." *Digital Journalism*, vol. 3, no. 3, 2015, pp. 398-415.

[16] McElreath, Richard. *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. Chapman and Hall/CRC, 2020.

[17] Parasuraman, Raja, and Dietrich H. Manzey. "Complacency and Bias in Human Use of Automation: An Attentional Integration." *Human Factors*, vol. 52, no. 3, 2010, pp. 381-410.

[18] Brooks, Frederick P., Jr. *The Mythical Man-Month: Essays on Software Engineering, Anniversary Edition*. Addison-Wesley Professional, 1995.

[19] Deming, W. Edwards. *Out of the Crisis*. MIT Press, 2000.

[20] Winner, Langdon. "Do Artifacts Have Politics?" *Daedalus*, vol. 109, no. 1, 1980, pp. 121–36.

[21] Trist, Eric L., and Ken W. Bamforth. "Some Social and Psychological Consequences of the Longwall Method of Coal-Getting." *Human Relations*, vol. 4, no. 1, 1951, pp. 3-38.

[22] Greenleaf, Robert K. *Servant Leadership: A Journey into the Nature of Legitimate Power and Greatness*. Paulist Press, 2002.

[23] Beck, Kent, et al. "Manifesto for Agile Software Development." *Agile Manifesto*, 2001, agilemanifesto.org.

[24] Brignull, Harry. "Dark Patterns: Deception vs. Honesty in UI Design." *A List Apart*, 2011.

[25] Tufekci, Zeynep. "YouTube, the Great Radicalizer." *The New York Times*, March 10, 2018.

[26] Rosenblat, Alex. *Uberland: How Algorithms Are Rewriting the Rules of Work*. University of California Press, 2018.

[27] Edmondson, Amy C. *The Fearless Organization: Creating Psychological Safety in the Workplace for Learning, Innovation, and Growth*. Wiley, 2018.

[28] Norman, Don. *The Design of Everyday Things: Revised and Expanded Edition*. Basic Books, 2013.

[29] Rother, Mike, and John Shook. *Learning to See: Value Stream Mapping to Add Value and Eliminate Muda*. Lean Enterprise Institute, 2009.

[30] Kant, Immanuel. *Groundwork of the Metaphysics of Morals*. 1785.

[31] Calvo, Rafael A., and Dorian Peters. *Positive Computing: Technology for Wellbeing and Human Potential*. MIT Press, 2014.

[32] Norman, Don. *The Design of Everyday Things: Revised and Expanded Edition*. Basic Books, 2013.

[33] Følstad, Asbjørn, and Petter Bae Brandtzæg. "Chatbots and the New World of HCI." *Interactions*, vol. 24, no. 4, 2017, pp. 38-42.

[34] Mori, Masahiro. "The Uncanny Valley." Translated by Karl F. MacDorman and Norri Kageki, *IEEE Robotics & Automation Magazine*, vol. 19, no. 2, 2012, pp. 98-100.

[35] Faden, Ruth R., and Tom L. Beauchamp. *A History and Theory of Informed Consent*. Oxford University Press, 1986.

[36] Schweitzer, Maurice E., and Rachel Croson. "Curtailing Deception: The Impact of Direct Questions on Lies and Omissions." *International Journal of Conflict Management*, vol. 10, no. 3, 1999, pp. 225-248.

[37] Weick, Karl E., and Kathleen M. Sutcliffe. *Managing the Unexpected: Resilient Performance in an Age of Uncertainty*. Jossey-Bass, 2007.

[38] Taleb, Nassim Nicholas. *Antifragile: Things That Gain from Disorder*. Random House, 2012.

[39] Sarter, Nadine B., David D. Woods, and Charles E. Billings. "Automation Surprises." *Handbook of Human Factors and Ergonomics*, 2nd ed., edited by Gavriel Salvendy, Wiley, 1997, pp. 1926-1943.

[40] Raji, Inioluwa Deborah, et al. "Closing the AI Accountability Gap: Defining an End-to-End Framework for Internal Algorithmic Auditing." *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020, pp. 33-44.

[41] Basiri, Ali, et al. "Chaos Engineering." *IEEE Software*, vol. 33, no. 3, 2016, pp. 35-41.

[42] Kahneman, Daniel. *Thinking, Fast and Slow*. Farrar, Straus and Giroux, 2011.

[43] Chesterton, G. K. *The Thing: Why I Am a Catholic*. Dodd, Mead & Company, 1929.

[44] McElreath, Richard. *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. Chapman and Hall/CRC, 2020.

[45] Weick, Karl E. *Sensemaking in Organizations*. SAGE Publications, 1995.

[46] Klein, Gary. *Sources of Power: How People Make Decisions*. MIT Press, 1998.

[47] Tufte, Edward R. *The Visual Display of Quantitative Information*. Graphics Press, 2001.

[48] Senge, Peter M. *The Fifth Discipline: The Art & Practice of The Learning Organization*. Doubleday/Currency, 2006.

[49] Christensen, Clayton M., et al. "Know Your Customers' 'Jobs to Be Done'." *Harvard Business Review*, September 2016.

[50] Brown, Tim. "Design Thinking." *Harvard Business Review*, June 2008.

[51] Liker, Jeffrey K. *The Toyota Way: 14 Management Principles from the World's Greatest Manufacturer*. McGraw-Hill, 2004.

[52] Merleau-Ponty, Maurice. *Phenomenology of Perception*. Translated by Colin Smith, Routledge, 1962.

[53] von Bertalanffy, Ludwig. *General System Theory: Foundations, Development, Applications*. George Braziller, 1968.

[54] Beck, Kent, et al. "Manifesto for Agile Software Development." *Agile Manifesto*, 2001, agilemanifesto.org.

[55] Christensen, Clayton M., et al. "Know Your Customers' 'Jobs to Be Done'." *Harvard Business Review*, September 2016.

[56] Fetterman, David M. *Ethnography: Step-by-Step*. SAGE Publications, 2010.

[57] Senge, Peter M. *The Fifth Discipline: The Art & Practice of The Learning Organization*. Doubleday/Currency, 2006.

[58] Liker, Jeffrey K. *The Toyota Way: 14 Management Principles from the World's Greatest Manufacturer*. McGraw-Hill, 2004.

[59] Rother, Mike, and John Shook. *Learning to See: Value Stream Mapping to Add Value and Eliminate Muda*. Lean Enterprise Institute, 2009.

[60] Based on Dave Thomas feedback session, 2024. Thomas warned against prescribing visual representations specifically, suggesting instead that the valuable representation is "whatever hurts most to produce."

[61] Sweller, John. "Cognitive Load Theory, Learning Difficulty, and Instructional Design." *Learning and Instruction*, vol. 4, no. 4, 1994, pp. 295-312.

[62] Tufte, Edward R. *The Visual Display of Quantitative Information*. Graphics Press, 2001.

[63] Brown, Tim. "Design Thinking." *Harvard Business Review*, June 2008.

[64] Star, Susan Leigh, and James R. Griesemer. "Institutional Ecology, 'Translations' and Boundary Objects: Amateurs and Professionals in Berkeley's Museum of Vertebrate Zoology, 1907-39." *Social Studies of Science*, vol. 19, no. 3, 1989, pp. 387-420.

[65] Brooks, Frederick P., Jr. "No Silver Bullet: Essence and Accidents of Software Engineering." *Computer*, vol. 20, no. 4, 1987, pp. 10-19.

[66] Beck, Kent, et al. "Manifesto for Agile Software Development." *Agile Manifesto*, 2001, agilemanifesto.org.

[67] Senge, Peter M. *The Fifth Discipline: The Art & Practice of The Learning Organization*. Doubleday/Currency, 2006.

[68] Beck, Kent, et al. "Manifesto for Agile Software Development." *Agile Manifesto*, 2001, agilemanifesto.org.

[69] Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.

[70] Hao, Karen. "Training a single AI model can emit as much carbon as five cars in their lifetimes." *MIT Technology Review*, June 6, 2019.

[71] Liker, Jeffrey K. *The Toyota Way: 14 Management Principles from the World's Greatest Manufacturer*. McGraw-Hill, 2004.

[72] Strubell, Emma, Ananya Ganesh, and Andrew McCallum. "Energy and Policy Considerations for Deep Learning in NLP." *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 3645-3650.

[73] Goldratt, Eliyahu M., and Jeff Cox. *The Goal: A Process of Ongoing Improvement*. North River Press, 2014.

[74] Schwartz, Roy, et al. "Green AI." *Communications of the ACM*, vol. 63, no. 12, 2020, pp. 54-63.

[75] The Standish Group. *CHAOS Report*. 2015.

[76] Japan Management Association. *Kansei Engineering: Understanding Consumer Feelings and Sentiment*. CRC Press, 2016. (This source discusses the broader Japanese philosophy of craftsmanship and quality, including Monozukuri principles.)

[77] Taleb, Nassim Nicholas. *Antifragile: Things That Gain from Disorder*. Random House, 2012.

[78] Ries, Eric. *The Lean Startup: How Today's Entrepreneurs Use Continuous Innovation to Create Radically Successful Businesses*. Crown Business, 2011.

[79] Kotter, John P. *Leading Change*. Harvard Business School Press, 1996.

[80] Chesterton, G. K. *The Thing: Why I Am a Catholic*. Dodd, Mead & Company, 1929.

[81] Baldwin, Carliss Y., and Kim B. Clark. *Design Rules: The Power of Modularity*. MIT Press, 2000.