



An Executive Field Guide to Implementing Trustworthy AI.

Information → Influence → Intention (Triple III Model) Al must earn trust before it earns autonomy By Dr. James L. Norrie, DPM, LL.M | October 13, 2025 Series Overview

We previously outlined our SAFER AI framework for ethical and responsible systems design. This series moves from principle to practice: how to make AI trustworthy in the real world and how to convert good intent into worthy action at the moment of human decision. A follow-on series will detail strategic use cases drawn from pilot results.

Experience suggests most enterprise teams try to "fix" AI in the wrong ways, something we see all the time. Teams chase better content—more data, faster models, slicker workflows—then wonder why people still do not rely on the system when it counts most. That is psychologically ignorant of basic scientific facts about how humans gather information and make decisions. Our view is therefore very different: higher human trust in AI is the missing ingredient. To close that gap, we focus on context—who the person is, how they decide, and what it takes for an agent to earn trust before exercising any autonomy, especially when stakes are high. When content and context intersect at the point of guery and reply, we proved AI output becomes human outcomes.

Why This Works (and Why Others Stall)

In our view, human delegation to agentic AI cannot be assumed. Or forced as a default setting because trust rapidly erodes as users detect unexpected gaps. Instead, autonomy in our model is earned through five specific, sequential trust gates— Provenance, Fit (myQ®), Stakes, Reversibility, and Ethical Alignment—which determine if and how an agent should act. Clear alignment with high confidence permits limited autonomy. Ambiguous alignment downgrades autonomy and requires explicit confirmation. Failed alignment declines the move and explains why. At every rung the human can slow, pause, or roll back. If any gate fails, the system reverts to counsel, not control. Autonomy becomes a privilege earned by performance under constraint.

Most solutions optimize what a model perfunctorily says. Triple III optimizes when to say it, how to say it, and whether to act at all. That is why it improves compliance, reduces accidental-insider risk, and raises adherence in finance, healthcare, education, and other domains where trust and influence are inseparable and outcomes are measured in avoided errors and on-time follow-through. Early pilots show that when trust and fit rise together, people are more willing to consider, comply with, and continue safer, higherquality actions.







Understanding the Triple III Model

Our patented Triple III Model operationalizes this shift through a staged, testable flow: Information \rightarrow Influence \rightarrow Intention. First, the system earns trust with verifiable evidence. Only then does it adapt how it communicates to fit the person in front of it. Finally, it converts guidance into a safeguarded plan with preview, confirm, and undo. Gates sit between steps. If trust falls below threshold, persuasion pauses, and the agent returns to evidence and context assessment. Influence stops feeling like pressure and starts functioning as collaboration.

Step 1: Improve Information Reliability (earn trust)

Harden the evidence tier with inspectable provenance, retrieval over a curated corpus, and honest calibration in place of confident guesswork. When sources disagree, defer to the system of record and show your work. The result is fewer hallucinations, clearer uncertainty, and answers that withstand scrutiny. If this gate does not clear, do not persuade; return to the evidence.

Step 2: Influence Through Style Alignment (earn a hearing)

Advice is accepted when it fits how people actually think. Our patented myQ® framework models durable differences in how individuals weigh risk, rules, and rewards, then maps those traits to reply style: tone, framing, evidence density, autonomy level, and challenge intensity. Same facts, different on-ramps. A rules-oriented user sees policy cites and checkpoints; a high-reward user gets the payoff and a clean path; a low-risk user sees limits, preview, and undo. This is personality, not persona theater—measurable, ethical, and programmable.

Step 3: Convert Persuasion to Intention (earn accountable action)

Once trust and fit are established, the agent helps users commit to a concrete plan matched to reversibility and stakes. This improves the likelihood of voluntary behavior change. Still at high-risk moments add friction through counterarguments, dual-source verification, and time-boxed holds. Low-risk, reversible steps move faster. Everything is logged in human-readable form, so accountability is visible to the user and auditable. Basically, autonomy is earned, never assumed.

Why Act Now?

As organizations scale agentic AI, users already sense which systems are generic talkers and which are dependable collaborative partners. This trust gap will widen, and trust is the entry ticket to influence. Teams that treat trust as the prerequisite, personality fit as the amplifier, and intention as the conversion step will set the new standard. They will





improve compliance, reduce risk, and lift human outcomes when it matters most, creating brand and economic advantages for those who lean in. If you are piloting or deploying AI and want to translate this method into your domain, let's talk. The clock is ticking, and the window to set the bar higher is open now.

First Up in the Series:

Before AI can persuade, it must prove it deserves to be heard. This opening article reveals why mistrust in chatbots is rational and how real trust begins with verified evidence, calibrated confidence, and psychological fit, not flattery. You'll get a first look at our patented myQ® framework, which links human decision styles with AI reliability. If you want AI that earns confidence instead of demanding it, this is where the blueprint begins.

Author Bio:

Dr. James L. Norrie is a professor of Law and Cybersecurity, and Founding Dean of the Graham School of Business, at York College of Pennsylvania (http://www.ycp.edu). He is the author of Beyond the Code: Al's Promise, Peril and Possibility for Humanity (Kendall, Hunt 2025). Learn more about our free community of interest in ethical Al at: www.techellect.com or visit www.cyberconlQ.com to learn more about Al tools to keep your employees safer online. To purchase his book, click on the QR code, or visit: https://he.kendallhunt.com/product/beyond-code-ais-promise-peril-and-possibility-humanity









Trust isn't a Feature; it's the Entire Game.

Information \rightarrow Influence \rightarrow Intention (Triple III Model) Al must earn trust before it earns autonomy By Dr. James L. Norrie, DPM, LL.M | October 13, 2025 Article I

The public remains uneasy with chatbots, and with good reason. Too often they are generic, error-prone—sometimes to the point of outright hallucination—and offer mostly trite, sycophantic replies that trigger derision. That lack of authenticity is baked in: they are trained to mimic us, but they are not us, lacking self-awareness. That machine hubris breeds mistrust and often kills influence at the very moment when influence could matter most, especially when the stakes of human error are high. My view is straightforward: Al should be an adjunct to human intelligence, not its replacement. Collaboration requires trust, and trust lives deeper than simple information retrieval. It begins with verified facts, but it rises or falls on whether advice fits how we actually think, argue, and decide. To gain influence with humans, AI must vastly improve its ability to earn our trust first.

The Science of Human Influence

Our solution to this is research-based. It blends psychology and technology in ways that move well beyond cosmetic "tone knobs." On the psychology side, my patented myQ® assessment maps how people navigate risk, rules, and rewards—the very essence of human decision-making. On the technology side, we harden the information layer with provenance and calibration, using retrieval from curated data sources to anchor reliability. Only then does style alignment follow. This way, the AI agent earns the right to advise by proving accuracy first, then shaping its reply to fit the user's decision tendencies. Not flattery. Not stereotypes, but a disciplined personal fit.

Additionally, if the stakes rise, the agent slows, showing more sources, and presenting counter-arguments to temper snap judgments driven by urgency. If the action is reversible, it lowers the safety net and speeds up action. That is how influence works without manipulation, and in harmony with how humans, not machines, approach decisions.

Why does this matter? Because mistrust in AI is rational. We should not hand judgment to machines that guess, pander, or bury their uncertainty. Even unintentionally, this can Iull humans into unwarranted false confidence. Yet there are moments when a persuasive, trustworthy agent can prevent real harm: a fraudulent transfer stopped in its tracks, a risky click avoided, a treatment finally taken on time. In practice, style-aligned agents built on verified evidence can improve compliance and reduce "accidental







insider" risk, without shaming users who might hesitate to ask human colleagues for help. That same pattern travels easily to healthcare, education, and government—domains where trust and influence are joined at the hip.

But let's be clear. The very approach that builds trust can also be weaponized. Used deliberately, it becomes deception: an agent tuned to overcome human defenses and extract trust where none is warranted. That is a risk we must control. Responsible Al should not only avoid these tactics but help train human users to recognize them, by flagging when an opposing agent violates ethical guardrails. The same levers that keep us safer can also be turned against us, and it is our responsibility to ensure they are not.

A One-Line Mathematical Model for Improving Al Trust

Here is a preview of our approach, a plain-language window into the complex engine underneath our efforts to make AI more trustworthy. The patent-pending details are more complex than what I can share here, but the intention can be captured in a simplified single line equation:

Reliance Score = Trust + Fit + Situation - Red Flags (all normalized to a 0–1 scale)

The specifics—variables, ranges, weights, thresholds—remain proprietary for now, but the general contours are easy enough to grasp:

- Trust: the quality of evidence, honestly adjusted for uncertainty.
- Fit: agentic alignment to style, tone, and framing with your risk, rule, and reward profile.
- Situation: the real-world context of stakes, time pressure, reversibility, and oversight.
- Red Flags: signals of weak provenance, prior errors, deception signals, policy violations.

If trust falls below a threshold, the agent reduces persuasion efforts. It returns to the evidence, flags the problem, and re-engages the human. Only after information trust is earned does it adapt style to increase influence through fit. This improves the likelihood of you committing to a plan, still with preview, confirm, and undo options, that turns collaboration into accountable action.

An Imperfect Conclusion

This approach is not a promise of perfect or permanently trustworthy AI. It is a disciplined way to make agentic AI worthy of limited autonomy, if the human in the loop chooses to grant it. Consent for style profiling is explicit and always easy to revoke. Influence is capped on irreversible or high-risk moves. And one requirement remains constant: WWHD—What Would a Human Do? If a reasonable human would reject the nudge, so should the agent. We must hard-code this principle into the math so that machine learning exhibits restraint, not hubris.





In Beyond the Code, I argued that AI should amplify our judgment, not replace it. AI and IQ must collaborate aligned to human values and ethics. Across this series, I will unpack the reliance equation in real-world scenarios—from phishing defense to medication adherence—to show the measurable gains when trust, fit, and safeguards all work together to achieve some of that potential. But the agentic AI clock is ticking. Now is the time to build agents that earn trust one transparent interaction at a time, before we delegate human authority to them on our behalf.

Up Next in the Series:

What "Better Information" Really Means. Before AI earns the right to persuade, it must first prove its evidence can be trusted. We'll unpack the first gate in the III Model—how to make data visibly reliable through verified provenance, precise retrieval, calibrated confidence, and independent checks. You'll see how these four principles turn abstract "trust" into measurable system behavior, why the best AI agents sometimes refuse to answer, and how trustworthy information becomes the foundation for real human allegiance.

Author Bio:

Dr. James L. Norrie is a professor of Law and Cybersecurity, and Founding Dean of the Graham School of Business, at York College of Pennsylvania (http://www.ycp.edu). He is the author of Beyond the Code: AI's Promise, Peril and Possibility for Humanity (Kendall, Hunt 2025). Learn more about our free community of interest in ethical AI at: www.techellect.com or visit www.cyberconlQ.com to learn more about AI tools to keep your employees safer online. To purchase his book, click on the QR code, or visit: https://he.kendallhunt.com/product/beyond-code-ais-promise-peril-and-possibilityhumanity









From Answers to Allegiance: Why Trust Begins with Better Information.

Information → Influence → Intention (Triple III Model) Al must earn trust before it earns autonomy By Dr. James L. Norrie, DPM, LL.M | October 13, 2025 Article II

We cannot build human allegiance with machine noise. It must be built on visible evidence humans can lean on. In my III Model, influence does not even get a turn until information quality clears a trust gate. So, this article zooms in on that first gate, showing how to improve your data game: provenance, calibration, retrieval quality, and honest uncertainty. If you want AI that persuades when it matters, it all starts by making it trustworthy to a fault.

Trust Gate 1: What "Better Information" Really Means

The rule of this trust gate is simple: gather and present only verified evidence and curated recommendations. At the system level, compute T, the trust credit. If T falls below threshold T₀, stop. Do not persuade and return to the evidence gathering cycle. In Beyond the Code, I argued that AI should amplify human judgment, not replace it. Fulfilling that promise starts with defensible information sources that earn allegiance to real human values, not an AI replication.

Most of our client organizations have too much data, not too little. This volume problem triggers a coherence challenge, making the work required to curate AI trust using only qualified data more intricate and complex to execute. To help, think about this in four ordered parts that are testable, implementable, and visible to the human-in-the-loop user of AI. They are our audience.

- 1. Provenance you can verify. Tie every claim to a source the user can inspect. Prefer content with credentials or signed attestations. Record lineage, including origin, last update, and ownership. Establish a single source of truth for key data and when systems disagree, defer to data gleamed from the system of record, not the most convenient copy.
- 2. Retrieval that deserves the name. Use Retrieval-Augmented Generation over a curated corpus. Keep the index clean with deduplication, semantic chunking, and tags for time, authority, and sensitivity. Make retrieval time-aware. Favor freshness for reversible decisions and authority for irreversible ones. Tune recall and precision to the task.
- 3. Calibration, not confidence theater. Replace hedged prose with explicit probabilities. Track calibration error and show reliability curves. Allow abstention by replying "I don't know," paired with a clear path to some evidence as this builds more trust than overconfident guesses that collapse under human scrutiny.





4. Independent checks where harm is plausible. For flagged or high-stakes tasks, cross-check answers across models using rule-based validators. Test for vulnerable populations that may be more susceptible to undue influence and potential Al liability. If results diverge, show the delta and pause persuasion for human review.

A Compact Blueprint for Immediate Information Quality

You can only earn trust by operationalizing it. The four principles above become credible only when they are visible in day-to-day system behavior, not trapped in a policy deck. This data trust blueprint turns those ideals into repeatable practices the human in the loop can see and verify.

- Corpus hygiene. Curate, tag for authority and freshness, and expire content on a schedule. Chunk on semantic boundaries so retrieval returns coherent ideas, not fragments. Carefully apply single source of truth principles throughout the corpus.
- Retrieval policy. Raise the bar when stakes rise. Require two independent sources or one authoritative source plus a validator for high-risk moves. Ask clarifying questions when intent is ambiguous to improve human perception of data relevance to query.
- Response construction. Lead with the answer and cite sources. State uncertainty and what would change the conclusion. Offer one safer alternative by default. Provide a "verify this" control that re-runs retrieval with stricter filters.
- Privacy and access controls. Enforce permissions before retrieval, not after generation. Redact at the source. Log queries, consulted sources, feedback and returns for data audit.

Get these four tactics right and Trises predictably, which means influence stops feeling like pressure and starts feeling like a trust-based collaboration between the human and the Al agent.

Two Brief Vignettes

Psychologically, humans experience truth and tone together, not sequentially as the demands of machine learning models require when deploying Al platforms. If the evidence tier is sloppy, the style tier reads as pandering. Existing systems that ignore this are already primed to be sycophantic, as I detailed in prior articles. Get the information right and you earn the right to speak in a way people can hear. That is why **T** unlocks **F** in our staged influence model.

Finance operations, late afternoon. A "rush" wire request appears. Retrieval pulls vendor history, contract terms, and the last approved routing data. Provenance is mixed. Calibration drops. $T < T_0$. The agent refuses to persuade, proposes a two-step fraud check, and places a time-boxed hold. A minute later the real CFO messages: do not wire.





Clinician support during rounds. A dosing question hits the agent. Retrieval anchors to the hospital formulary and current guidelines. Calibration is high and recent. $T \ge T_0$. The agent answers plainly, cites both sources, and offers a dosage calculator. It shows a reversible path: "Place order with pharmacist review."

In both vignettes, the construction of information gate did the real work. The agent earned influence by being predictably right or predictably cautious, and that consistency compounds into higher human trust. And this model travels—finance, healthcare, education, government, and beyond. This is not a generic one size fits all agentic AI platform; rather, it is a code-ready trust and influence methodology that can power your own AI pilots and use cases to successful implementation and results by improving user trust and measurable influence over time.

Tracking Metrics That Actually Move Trust

You cannot manage what you will not measure. In the data-centric AI era, the issue is rarely inability to measure. It has been reluctance. Accumulating and aggregating information was seen as progress. That mindset has left many organizations data-saturated and wisdom-deprived. Multiple, less reliable sources spur scattered employee and customer actions, influenced by data designed more for manipulation or convergence on a viewpoint than for trust and influence. To fix this, strengthen data management with a short, credible list:

- Source validity rate. Percent of citations that pass authority and freshness checks.
- Provenance completeness. Lineage fields present and accurate per answer.
- Calibration error. Gap between predicted and observed correctness over time.
- Model-agreement delta. Disagreement rate on high-stakes prompts, routed to review.
- First-error recovery time. Speed to detect, correct, and notify after a bad answer.

Organizations that invest here deliberately advance data quality and improve their appropriate reliance ratio. This is a strategic measure of how often users flag responses as unreliable for any reason. That ratio should decline as quality rises, with lower being better. It means humans lean on the agent because it is right more often than it is wrong, and on-going allegiance is earned.

Bringing it Back to the III Model

Our model's decision flow is not personality theater. It is a data brake pedal. Only with data reliability can we expect sustained, measurable style-aligned influence to emerge.

- 1. Information. Verify evidence and compute T.
- **2.** Gate. If $T < T_0$, return to Information. Do not persuade.
- **3. Influence.** Shape style to the person and context. Compute **F.** Trigger a challenge protocol on high-risk or irreversible decisions.





- **4.** Intention. Turn guidance into a concrete, safeguarded plan of action.
- 5. Act and log. Execute with preview and undo. Update T and F from outcomes.

Up Next in the Series:

We move to the second stage, Influence, and show how style alignment works without flattery. We will translate risk, rules, and rewards into practical reply strategies, and explain the challenge protocol that prevents "influence" from becoming manipulation. The clock is still ticking. Let's keep earning trust, one transparent interaction at a time.

Author Bio:

Dr. James L. Norrie is a professor of Law and Cybersecurity, and Founding Dean of the Graham School of Business, at York College of Pennsylvania (http://www.ycp.edu). He is the author of Beyond the Code: AI's Promise, Peril and Possibility for Humanity (Kendall, Hunt 2025). Learn more about our free community of interest in ethical AI at: www.techellect.com or visit www.cyberconlQ.com to learn more about AI tools to keep your employees safer online. To purchase his book, click on the QR code, or visit: https://he.kendallhunt.com/product/beyond-code-ais-promise-peril-and-possibilityhumanity









Stop Persona and Start Personality.

Information \rightarrow Influence \rightarrow Intention (Triple III Model) Al must earn trust before it earns autonomy By Dr. James L. Norrie, DPM, LL.M | October 13, 2025 Article III

We do not trust cardboard cutouts; we trust people. If agentic AI expects a hearing, it must stop performing as a generic persona and align to real human personalities. Not sycophantic flattery and all-too-easy stereotypes. Rather, a disciplined fit with how individuals actually think, argue, and decide things for themselves even when collaborating with machines.

In my research for Beyond the Code, I asked how AI might amplify human judgment rather than replace it. That became a central thesis of the book. Now refined into the "Triple III Model," the work shows in practice how curated information earns trust, how influence increases when aligned to personality theory, and how that, in turn, raises intention to convert sound guidance into safeguarded action. It changes decision behavior. This article focuses on that second step—influence—and the single insight that unlocks it: deeper personalization increases influence only when it is grounded in measurable human personality traits, not in superficial machine personas.

Why Personas Fail & Personality Works

Personality is not a mood or a marketing segment. It is a durable, measurable, predictive pattern of traits that shows up in how people weigh risk, respect rules, and respond to rewards. That triad, modeled through the myQ® framework, explains why some of us demand authoritative policy cites and checkpoints before acting, why others pursue a clear payoff path that balances outsized risk for outsized reward, and why still others insist on a visible "undo" to test a choice's intuitive feel before committing. Advice that meets those tendencies will be heard; advice that pushes against this intrinsic human architecture will be ignored, no matter how clever the prose. Replicating this complexity in an AI platform is possible, but it requires a deep, working grasp of both psychology and AI technologies.

By contrast, the slick personas that dominate today's agentic AI experiments are efficient fictions at best and superficial sycophancy at worst. They compress people into tidy labels that may suit a campaign but fail at moments of consequence. A persona can pick a color palette; it cannot support a chemotherapy decision, approve a wire transfer with accountability, or help a teacher intervene at the right moment. That requires deeper alignment to personality, which the myQ® framework provides by giving AI agents a scientifically grounded map to align with user style.







From Theory to Practice: Programming Personality without the Gimmicks

Patents are scrutinized for novelty and substance; they are not awarded for "vibe dials." Dismissing the myQ® framework as cosmetic entirely misses the point. Instead, focus on its two proven pillars of psychology theory:

Trait-based personality science. Decades of research show that stable traits shape how people process information, evaluate risk, and follow through. That yields a credible, durable map of meaningful interactional differences measurable across time, content and context.

Cognitive bias and decision psychology. The same levers that can protect us—urgency, authority, social proof, scarcity—are often exploited to hack judgment. We codify those levers, use them transparently, and hold the model to account when deploying it.

For machines to approximate human personality usefully, alignment must be programmable, testable, and improvable in the digital wild. In the III Model, information still comes first; fit must never outrun truth or trust collapses. Once the trust gate clears, the next gate is influence—the agent's ability to speak so its human collaborator can actually hear it, empathetically and effectively. Here is the practical training sequence:

- 1. Consent or a clear cold start. The user opts into profiling via a validated myQ® assessment linked to their Al profile, or the agent begins from a clearly labeled cold start, using sparse, provisional signals that improve with use.
- 2. Map personality to a reply profile. The user's myQ® vector across risk, rules, and rewards sets tone, framing, evidence density, autonomy level, and challenge intensity—reframing basic information into responses the user experiences as more influential.
- **3. Keep trust ahead of fit.** If the information tier drops below threshold during an interaction, persuasion pauses. The agent returns to evidence, validates sources, and states uncertainty plainly.
- **4. Validate fit as a living hypothesis.** The goal is improved comprehension and follow-through for this person in this context. When alignment doesn't help, the agent adjusts and re-tests.

What follows is a cycle that closes the loop with the decision-maker, now feeling more helpful to the user rather than potentially intrusive or pushy. Instead of locking people into static labels, the agent watches how alignment performs over time and adapts accordingly.

Can you imagine the difference? A generic agent talks **at** you; a personality-attuned agent works **with** you. Plans hold without late reversals because guidance fits how you weigh risk and stays within your tolerances. When a counter-argument surfaces in a





high-stakes moment, the agent slows the decision just enough to reveal the safer path. Confidence rises for the right reasons—clear sources, reversibility, and preview-andundo—rather than because the prose sounds certain. And when the system errs and owns it, trust recovers.

These are not vanity signals or easy agreement. They are the recurring feedback that steers your reply profile toward your stable pattern, creating a self-training, continuously improving fit anchored to evidence, not applause. We keep one plain test in view: WWHD—What Would a Human Do? Personality alignment should help people act wisely, not merely agree more often. If a reasonable human would reject the nudge, so should the agent. That builds and keeps trust!

Your Monday Morning Moves

Retire the persona approach in applications intended for human influence or decision support. Instead, stand up a personality-aware pilots with our myQ® framework, using opt-in and a visible transparency with users to improve trust and influence. Wire in counterarguments for high-stakes moments and track what matters: comprehension, safe-path selection, adherence without late overrides, and trust recovery after early errors. Share those results with users and market the advantages of style-aligned personalized agents versus generic Al. When people can see you earning their trust and see the system adjusting to them—they will lean into its value. And you will reap the business value you always new AI could have but hasn't yet achieved.

Summary Conclusion

Think about the moments when judgment wobbles; a late night, a crowded inbox, a decision that matters requiring more than the clock allows. What can help in those moments? An Al-enabled tool that knows not only the relevant facts but inherent knowledge of the person weighing them. This approach truly demonstrates human care and concern expressed in collaborative AI tools.

Personality is the map; but better human decision-making is the journey. When Al honors both, particularly matching the way you balance risk, rules, and rewards to what is actually true, its counsel stops feeling like pressure and starts reading as useful partnership. That is not a tone trick or a superficial friendly veneer. It is the quiet architecture of collaboration: guidance that meets you where you are and walks with you just far enough to make the hard thing easier to do. If you are building AI, the clock is ticking; treat personality alignment as a critical best-in-class goal, and if you want help turning that principle into practice, we can show you how.





Up Next in the Series:

The final article in the Triple III series explains when AI truly deserves the right to act, and when it must step back, by arguing that autonomy must be earned. You'll get a tight, practical ladder of safeguards, plus real-world rules showing how agents should ask for handoffs, offer undo paths, and always let people say "no" without penalty. If you care about AI that helps rather than hurries, this article shows the code-ready guardrails to make delegation safe, auditable, and human-centered.

Author Bio:

Dr. James L. Norrie is a professor of Law and Cybersecurity, and Founding Dean of the Graham School of Business, at York College of Pennsylvania (http://www.ycp.edu). He is the author of Beyond the Code: AI's Promise, Peril and Possibility for Humanity (Kendall, Hunt 2025). Learn more about our free community of interest in ethical AI at: www.techellect.com or visit www.cyberconlQ.com to learn more about AI tools to keep your employees safer online. To purchase his book, click on the QR code, or visit: https://he.kendallhunt.com/product/beyond-code-ais-promise-peril-and-possibilityhumanity









Guardrails to Autonomy: When Humans Delegate and When They Decline.

Information → Influence → Intention (Triple III Model) Al must earn trust before it earns autonomy By Dr. James L. Norrie, DPM, LL.M | October 13, 2025 Article IV

Our article series has traced the link between emerging technology and established psychology—AI and IQ working together as AIQ, particularly in decision-support settings where influence on human decision-making is the goal. Collaboration with AI is the promise, but decision delegation is the test. Any AI system should encourage handoff only when people are comfortable with the potential consequences, because those consequences land on humans, not machines.

This final article in the "Triple III" model sets the conditions under which an agent earns limited autonomy and the safeguards that keep that autonomy answerable to human values. The question is straightforward and consequential: under what conditions should a person confidently hand off a decision to an agent, and when should they refuse? In Beyond the Code I argued for machine-human collaboration in which ethical human judgment stays in charge. Here we close the loop by specifying when that AI collaboration should, and should not, extend to decision delegation.

Autonomy Must Be Earned, Not Assumed

Autonomy is not a toggle in a settings panel. It is a progression that moves with evidence of curated information quality, personality fit, and nudging restraint. In our model, trust and influence arrive in order: information earns initial trust; fit earns a hearing with higher human impact; intention turns guidance into an action plan at the moment of decision. Only after all those stages succeed should any system ask for autonomy to act. Think of it as a short ladder.

At the first rung, the agent retrieves, cites, and explains its curated sources, keeping them transparent and traceable. Organizations with a solid foundation of curated data can build this trust corpus and apply it to agent design through retrieval-augmented techniques. For many, the critical task remains sorting, classifying, and turning existing data into reliable AI training material. But without that work, AI agents will be seen as less useful and less accountable.

The second rung engages the personality traits of users based on their inherent perception of risk, rules, and rewards using our patented myQ® tool. That structures Alagent replies in well-studied ways to offer guidance that gains credibility and influence





by aligning with measurable personality tendencies and preferences. We have seen this succeed in cybersecurity, where the tool reduces risky behavior in real-time.

The next rung is the impact tier. The agent must assess the potential effect of a decision and behave differently on low-risk, fully reversible steps with preview and undo than it does on medium- or higher-risk steps. Techniques include decision brakes such as dualsource verification, visible counter-arguments, and clarifying questions before execution. The top rung is not autonomy at all. It is supervised action on high-risk choices that explicitly require human confirmation. Moving up or down these rungs is earned by performance under constraint, not by charm or convenience. These tiers map neatly to enterprise risk practices and open a path to stronger regulatory compliance. They also apply to consumer contexts where trust precedes influence, improving adoption and closing rates when the fit is real.

The right to say no sits beside every rung. People must be able to decline, pause, or roll back without penalty or judgment, or trust will lapse. The agent's job is to show what will happen, how to undo it, and what evidence would merit a different recommendation next time. Control remains with the person who bears the consequence or enjoys the reward; the difference is how and when the nudge is applied, and how the outcome is reviewed.

Trust and influence are inseparable. Systems operate in code, but they must serve humans who do not. When people trust the information and experience a fit that respects how they weigh risk, rules, and rewards using our patented myQ® lens, they are more willing to consider, comply with, or conform to actions that improve outcomes. This is essential to improving corporate compliance, increasing willingness to trust as a proxy for increased adherence to desired behaviors. The result is outsized impact in highstakes, high-reward settings across industries.

The Five Delegation Gates

In our model, delegation is a sequence of gates. They open in order, or they do not open at all.

- 1. Provenance. Are the sources strong and current for this task, and can the human inspect them if needed? If not, stop and return to the evidence.
- 2. Fit. Is guidance framed to align with the person's stable tendencies identified by the myQ® framework, so they can better hear it and judge it? If alignment is weak, adjust the presentation, and confirm with the user that the style reflects their preferences accurately.
- 3. Stakes. What is the worst credible outcome, and who bears it? Higher stakes lower the ceiling on autonomy and shift the system from forward nudging to decision braking to slow instinctive errors.



- 4. Reversibility. If a step misfires, can we unwind it quickly and fully? If yes, autonomy can rise within that boundary; if not, require deliberate human confirmation before proceedina.
- 5. Ethical Alignment. Would a reasonable human accept this nudge in this context, given these stakes, and how confident is the system in that judgment? Clear alignment with high confidence supports action at the appropriate autonomy level. Plausible alignment with low confidence should downgrade autonomy and require explicit confirmation. If alignment fails because the action would undermine welfare, dignity, or informed choice, decline and explain why.

If any gate fails to meet its threshold, the system reverts to counsel, not control, and reduces pressure to act. To maintain trust, especially once earned, the agent should avoid obfuscation, refuse to guess, and be transparent about where it falls short because of deficient or low-quality information. It should also detect when a user is at risk or when a conversation turns toward self-harm, or conflicts with legal or regulatory boundaries, and withhold help that would enable a negative outcome. A gated system protects trust while still offering timely help.

A Healthcare Example, End-to-End

A regional hospital considers allowing its agent to auto-approve noncritical medication refills during morning rounds. The system runs the gates in order. Provenance: are formulary and guidelines current and cited in the record? Fit: will clinicians see the basis for the suggestion in a format that matches how they weigh risk and rules, with a pharmacist review visible by default? Stakes: a refill error could inconvenience but not endanger, so risk is limited. Reversibility: cancellation is available with one click, and alerts are sent immediately. Ethical alignment: would a reasonable human accept this nudge here and now, and is confidence high enough to proceed without delay? With five gates cleared, the hospital authorizes the limited handoff. If any gate fails on a specific case, the agent stops, explains why, and asks for a human decision.

Setting a New AI Collaborative Systems Benchmark

The Triple III model is simple to state and increasingly practical to implement as capabilities grow. Expectations are rising as users spot deficiencies and expect them to be fixed. If we expect humans and AI to collaborate successfully, we must expect to trust these systems. Otherwise, why would people default to them for decision or compliance support? The pathway is not cosmetic. It is science-based, researchinformed, human-tested, and the results show stronger acceptance by end users. With improved data transparency, and the patented myQ® framework deployed appropriately, trust is earned through evidence and personality fit, not flattery. That earned influence is converted into intention, then preserved over time by ensuring levels of autonomy are granted with continuous safeguards, not slogans. The order matters



because people matter. When systems respect that order, collaboration feels like help rather than pressure, and handoffs feel like judgment shared rather than judgment surrendered.

Build What Comes Next

If you want AI that helps rather than hurries, treat delegation by the human to the machine as a privilege that must be earned and kept. Treat these gates as real gates, not theater. Let people say no without penalty and make reversibility the default. Show your work in language a non-expert can understand. If you are building or piloting systems and want to translate this model into practice in your specific domain, we should talk. AI is moving fast, and the teams that align autonomy with human values will set the standard everyone else has to meet.

Author Bio:

Dr. James L. Norrie is a professor of Law and Cybersecurity, and Founding Dean of the Graham School of Business, at York College of Pennsylvania (http://www.ycp.edu). He is the author of Beyond the Code: Al's Promise, Peril and Possibility for Humanity (Kendall, Hunt 2025). Learn more about our free community of interest in ethical Al at: www.techellect.com or visit www.cyberconlQ.com to learn more about Al tools to keep your employees safer online. To purchase his book, click on the QR code, or visit: https://he.kendallhunt.com/product/beyond-code-ais-promise-peril-and-possibility-humanity



