



# **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 dual-source 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 high-stakes, 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 (<a href="http://www.ycp.edu">http://www.ycp.edu</a>). 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: <a href="https://www.techellect.com">www.techellect.com</a> or visit <a href="https://www.cyberconlQ.com">www.cyberconlQ.com</a> to learn more about Al tools to keep your employees safer online. To purchase his book, click on the QR code, or visit: <a href="https://he.kendallhunt.com/product/beyond-code-ais-promise-peril-and-possibility-humanity">https://he.kendallhunt.com/product/beyond-code-ais-promise-peril-and-possibility-humanity</a>



