



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 **T** rises predictably, which means influence stops feeling like pressure and starts feeling like a trust-based collaboration between the human and the AI 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 datasaturated 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



