



The Limits of Statistical-Learning AI Are Our Architecture

Autonomy Platform · Public Blog

CONFIDENTIAL

© 2026 Azirella Ltd. All rights reserved worldwide.

Strictly confidential and proprietary — do not distribute.

The Limits of Statistical-Learning AI Are Our Architecture

A widely-shared technical retrospective makes the honest case for what today's AI cannot do: out-of-distribution reliability, self-explanation, causal reasoning, and long-horizon decisions. In supply chain, those four limits are the spec we built to. Here is the mechanism for each, and the operating model that keeps a human in control.

A careful technical retrospective went around this month: [“The Second Wave of Artificial Intelligence: Inside the Statistical Learning Era That Built Everything You Now Call AI”](#) by Hayanan, on Medium. It walks through DARPA’s three-wave taxonomy of AI: handcrafted rules, then statistical learning, then a not-yet-here “contextual adaptation.” Its argument is that we are still deep inside the second wave, the statistical-learning era, and that the second wave is, at the bottom of the stack, one move repeated at ever-larger scale: minimize a loss function over data by gradient descent. The history is right. The honesty is rarer.

The part worth keeping is the asymmetry the author is willing to name. Statistical learning is superhuman at pattern recognition and sequence generation. It is structurally weak at four things:

1. **Out-of-distribution reliability** - performing well on inputs that differ from the training data.
2. **Explaining its own output** - producing a confident answer without being able to say why.
3. **Causal reasoning** - the difference between observing $P(y \text{ given } x)$ and predicting $P(y \text{ given } do(x))$, the effect of actually intervening.
4. **Long-horizon decisions** - carrying a multi-step plan to a reliably correct conclusion.

The field has real responses to each: retrieval grounding, mechanistic interpretability, neuro-symbolic methods, test-time compute that lets a model search and verify at inference. These are not nothing. But the retrospective's practical advice for everyone outside the frontier labs is narrow: be a consumer of frontier models, build applications on top. That advice misses a third category, and the third category is the whole point.

One architecture, two halves

Statistical learning is superhuman at four things and structurally weak at four others

Does well, often superhuman

- ✓ Pattern recognition
- ✓ Sequence generation
- ✓ Transfer learning
- ✓ Closed-world reinforcement learning

Structurally weak: the four blind spots

- ✗ Out-of-distribution reliability
- ✗ Explaining its own output
- ✗ Causal reasoning, $P(y \mid \text{do}(x))$
- ✗ Long-horizon decisions

Reworked from Hayanan's "two things the second wave can do and the two it cannot." The right-hand column is the specification the rest of this post answers.

The narrower, harder claim

A foundation model is one thing. An application sitting on top of one is another. There is a third thing: a **decision substrate** for a single domain, whose entire job is the four items on that list. We did not build it to answer a taxonomy. We built it because a planning agent that decides what to order, what to move, and what to make has to be correct for a reason, and a model running on a flat vector space cannot give you the reason.

Here is the mechanism for each of the four, in a form you could point at in the code.

Out-of-distribution reliability becomes a calibrated band. Every decision the substrate makes carries a conformal-prediction interval: evidence in, a calibrated probability out, with a coverage guarantee rather than a vibe. When the world drifts away from what the agent has seen, the

band widens and the system says so. That honesty is the input the next layer keys on. An agent that cannot tell you how sure it is has no business acting autonomously; an agent that reports a calibrated likelihood can be governed.

Causal reasoning becomes interventional data. The retrospective is exactly right that you cannot get $P(y \text{ given } do(x))$ from observation alone: you need experiments or structural assumptions. So we run the experiments. A digital twin of the supply network lets the agents act counterfactually millions of times - change the order, change the lane, change the buffer, run it forward, measure the outcome. That is reinforcement learning against a simulated world, and it is the engine that turns correlation into a policy that has been tested against intervention. A causal layer then attaches counterfactual effect estimates and override-effectiveness to real decisions. This is the piece that distinguishes a prediction from a decision.

Self-explanation becomes a contract. Every decision exposes four fixed fields on demand: what prompted it, what was decided, the expected outcome, and the calibrated likelihood that outcome is realized. Explanation is not a post-hoc story generated to please an auditor. It is the structured record the agent reasoned from, surfaced when a person asks.

Long-horizon decisions go to purpose-built, supply-chain-specific neural networks. The decision tier is not a general-purpose language model asked to plan. It is a family of narrow models trained, by reinforcement learning on the twin, to reason over the supply network itself: sites, products, lanes, partners, and the flows between them. The language model is confined to one job - narrating decisions that the decision tier has already made. It is never on the decision axis. Roughly

ninety percent of the substrate is decision machinery; the language layer is the last ten percent, and it talks about the work rather than doing it.

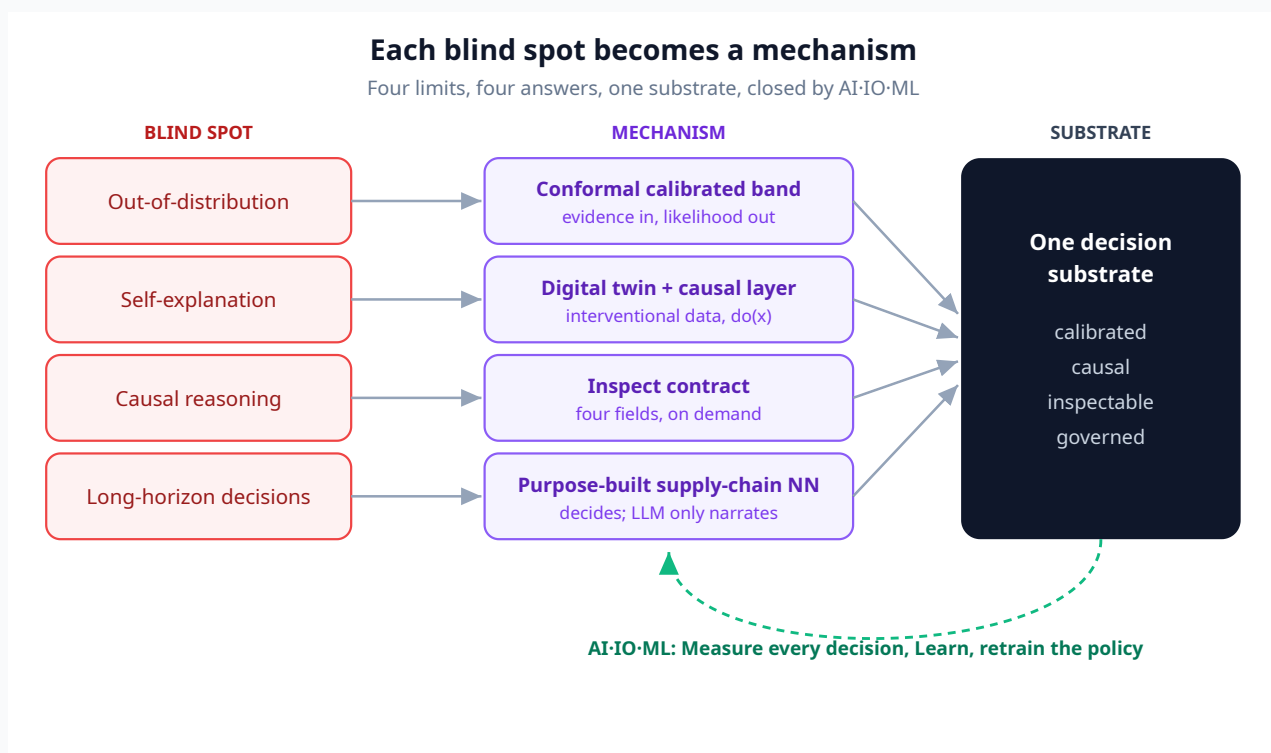
The operating model is what makes it safe to automate

A substrate that can act autonomously is only useful if a human stays in control without becoming the bottleneck, and only durable if it gets better with use. That is the [AI·IO·ML operating model](#): six verbs in three couplets.

- **AI - the agent acts: Automate, Inform.** The agent decides and acts within its declared envelope, with no approval queue, and informs a human when a decision crosses a policy boundary or when calibrated confidence is low and urgency is high. The conformal band from the first item above is exactly what the Inform threshold reads.
- **IO - the human engages: Inspect, Override.** A person can pull any decision and get the four-field explanation on demand, and can override it because they know something the agent did not. Override is not undo. It is a new, better-informed decision.
- **ML - the system improves: Measure, Learn.** Every decision and every override is measured against its outcome, and the result is learned: the policy retrains on the twin and the Inform threshold recalibrates, so calibration tightens and more decisions safely move into Automate every cycle.

The ML couplet is where the reinforcement-learning loop closes, and it is the part the retrospective's "consumer of frontier models" framing has no room for. An override is not just a correction to one plan; it re-enters the

loop as training signal. The system that decides and the system that learns are the same system, governed by one model. The human is out of the loop on routine execution and decisively in control of how the machine learns.



The four blind spots, each answered by a mechanism, composed into one decision substrate and closed by the AI·IO·ML loop. The dashed return path is the ML couplet: every decision and override retrains the policy.

Why the limitation list lines up with our pillars

This is not luck, and it is not retrofitting. We arrived at this shape from one rule we wrote years ago and have never relaxed: the language model is never the decider. Once you refuse to let a flat operations engine pretend to be the whole machine, everything else is forced. If the language model is not the decider, something has to be, and that something needs calibrated uncertainty to know when to ask, interventional data to reason

causally, an explanation contract to be inspectable, and a learning loop to improve. Those are not four features we collected. They are what you are obliged to build the moment you take the four limitations seriously instead of waiting for a larger training run to dissolve them.

Hayanan is honest that those limits are structural to statistical learning alone. I agree. The move is not to abandon statistical learning; it is to compose it with the three things it cannot supply by itself, inside a domain where being correct for a reason is the product. In supply chain, that substrate is running against real plans today.

The second wave's blind spots are not a warning for us. They are the specification.

See Autonomy in action

Walk through how Autonomy models, executes, monitors, and governs supply chain decisions with autonomous AI agents.

[See It Live](#)