

EON Switchboard

Best value money can buy. Frontier quality where it matters. Security you never compromise on the customer's behalf.

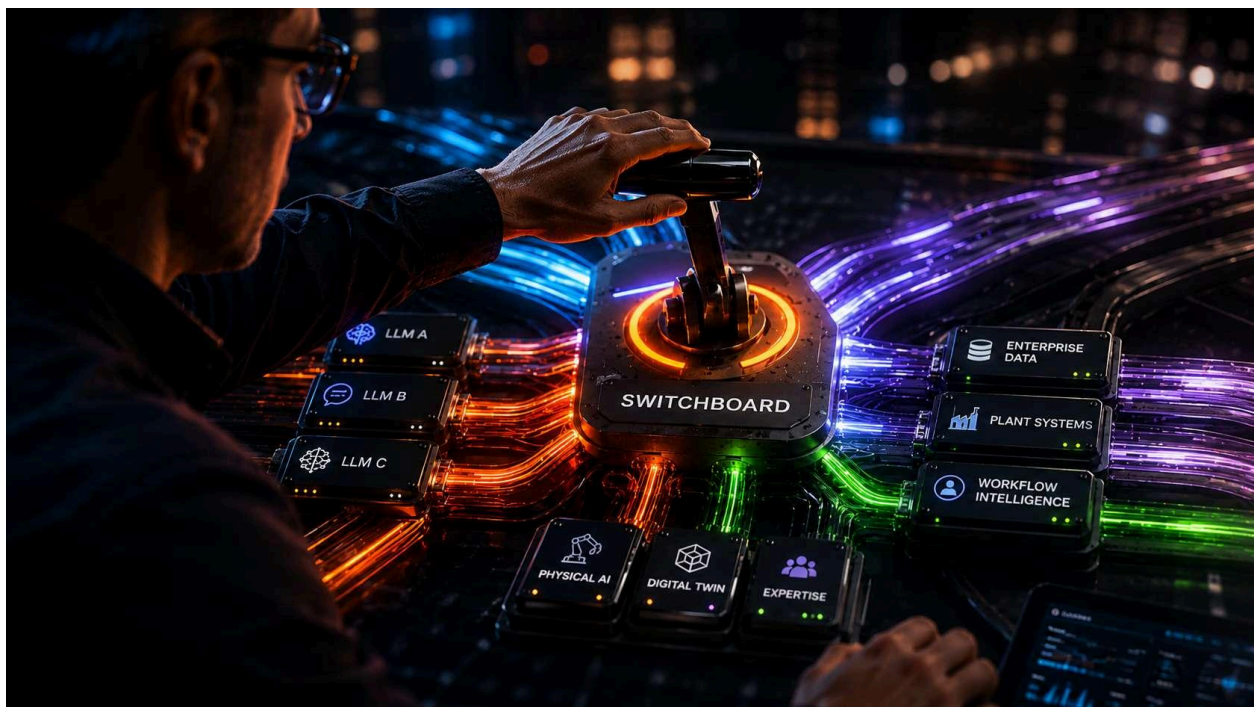


Table of Contents

- Executive summary..... 2**
- 1. The inflection: model commoditization and a slowing, hidden frontier..... 2**
- 2. Why one frontier model is the wrong architecture for Work Intelligence..... 3**
- 3. The EON Switchboard: three principles that never trade against each other..... 3**
 - Value..... 3**
 - Quality.....3
 - Safety.....4
- 4. Architecture: two planes..... 4**
 - The build plane — creating the products..... 4
 - The run plane — customer inference at scale..... 4
- 5. The decision flow..... 4**
- 6. The sovereignty gate: customer-controlled security posture..... 5**
- 7. The model bench, by job..... 5**
- 8. The economics of routing..... 6**
- 9. Why this is defensible: the proprietary-data moat..... 6**
- 10. Where Switchboard sits: Conductor, Switchboard, Verdict..... 6**
- 11. Conclusion..... 6**

Executive summary

Enterprise AI is at an inflection point. Frontier models are bunching together and, increasingly, are slow or restricted to release; open-weight models now perform within a hair of the frontier at a fraction of the cost; and intelligent routing between models has become standard infrastructure. EON Switchboard is the layer that turns this shift into customer value.

EON Switchboard is the model-routing and governance layer beneath every EON product. It receives each task, makes a fast judgment, and routes it to the right model under firm rules the customer sets — so customers receive the best value money can buy without ever trading away quality or safety. Its distinguishing feature is a customer-controlled sovereignty switch: enterprises choose, per tenant, which classes of model are permitted to run their work, and EON can prove what ran where.

This paper explains the market shift that makes Switchboard necessary, the two-plane architecture that separates building products from running them at scale, the decision flow and its sub-millisecond overhead, the sovereignty gate that sets EON apart, the economics of routing, and why this approach is defensible precisely because it is built for Work Intelligence and physical AI — problems that a single general-purpose LLM cannot solve on its own.

1. The inflection: model commoditization and a slowing, hidden frontier

Three things are true at once in mid-2026, and together they redraw the enterprise AI map.

First, the frontier is slowing and partly going dark. The most capable models are increasingly held back or restricted under export-control pressure; recent top-tier releases were suspended within days of launch, and at least one frontier release was held back from public availability. Enterprises can no longer assume the single best model will be available, stable, and permitted in every jurisdiction where they operate.

Second, open-weight models have nearly caught the frontier. On core benchmarks, the leading open-weight models now sit within roughly two-tenths of a point of the closed frontier, while costing as little as a thirtieth as much per token. Several are permissively licensed and can be self-hosted, which matters enormously for data sovereignty and air-gapped deployments.

Third, routing between models is now standard practice rather than a clever optimization. Teams that send each request to the cheapest model that can handle it report cost reductions in the range of forty to eighty-five percent with no visible drop in quality, because most production traffic never needed a frontier model to begin with. Industry analysts now describe AI gateways and routers as critical infrastructure rather than optional tooling, and the routing decision itself adds well under a millisecond when made with simple rules.

The combined message is clear: the model is becoming a commodity input. Competitive advantage is moving away from access to any single model and toward how intelligently an enterprise orchestrates many of them — and toward the proprietary data and work the models

are applied to.

2. Why one frontier model is the wrong architecture for Work Intelligence

The instinctive enterprise approach — pick the single best model and route everything to it — fails on three fronts, and each failure is amplified at Work-Intelligence scale, where millions of inferences run across many products.

- Expensive at scale. Paying frontier prices for the large majority of tasks that never needed a frontier model is pure waste, and that waste grows linearly with adoption — exactly as the value of AI to the customer grows.
- Brittle and single-vendor. One provider outage, one rate-limit, or one suddenly suspended model takes the whole product down. Recent suspensions show this is not hypothetical.
- A sovereignty and governance risk. Large enterprises will not accept dependency on models they cannot govern — for example, models whose provenance or licensing they have not approved — and they must be able to prove what model processed which data, in which jurisdiction.

Work Intelligence and physical AI raise the stakes further. Traditional LLMs answer questions; they do not, on their own, understand the work — the equipment, the procedure, the competence record. That understanding requires a composable ontology plus orchestration across many specialized models, gated for safety. The model is an input. The understanding is the product. A single general-purpose LLM cannot deliver it.

3. The EON Switchboard: three principles that never trade against each other

EON Switchboard is built so that value, quality, and safety are delivered together, not balanced against one another.

Value

Every task runs on the cheapest model that clears the quality bar for that specific job. Because the price spread between model tiers can be roughly thirty-fold, routing the right traffic to the right tier produces large, compounding savings on customer inference.

Quality

The work of building EON's products is always done on best-in-class frontier models; nothing customer-facing ships below a defined quality line, regardless of price. Cost optimization applies to running the products, never to compromising the output customers see.

Safety

A model's provenance and a customer's data posture are hard gates, evaluated before cost — never traded away for a cheaper route. This is the dimension most routers ignore, and it is where EON draws its line.

4. Architecture: two planes

The central design decision is to separate two fundamentally different workloads, because they have opposite economics.

The build plane — creating the products

Here EON always uses the best model money can buy. Build work is never cost-optimized, because the speed and quality of what EON ships compounds over time; a cheaper model on the build plane is false economy. Typical models: Opus 4.8, GPT-5.5, and Mythos-class models such as Fable 5 when available.

The run plane — customer inference at scale

Here the Switchboard earns its keep. Run-plane traffic is cost-optimized, compliance-gated, and subject to the customer's sovereignty switch. This is where value compounds: the cheapest qualified model handles each request, while the audit trail records exactly what ran. The guiding rule is simple — spend on the build, commoditize the run.

5. The decision flow

When a task arrives on the run plane, the Switchboard resolves it in four steps, with the routing decision itself adding well under a millisecond.

- Classify fast. A rule-first classifier reads task type and length — no additional model call — keeping overhead negligible.
- Apply the sovereignty gate, before cost. Anything the customer has not approved is denied by default. Safety and provenance are decided before price ever enters the calculation.
- Route on cost-for-quality. Within what the gate allows, the cheapest model that clears the bar for that job is selected.
- Fall back. If a chosen model is unavailable, the Switchboard fails over automatically. Routing and failover are kept distinct — EON does both.

Every routing decision is logged and auditable. The audit trail is not an afterthought; it is the substance of the security story EON can tell a regulated customer.

6. The sovereignty gate: customer-controlled security posture

The sovereignty gate is the capability that sets EON Switchboard apart. Each customer tenant is assigned a posture that governs which classes of model may run its work:

- US-only — American models exclusively; no foreign-provenance or export-restricted model touches the tenant.
- Allied / open-weight — US plus allied and permissively-licensed open-weight models.
- Global-South permissive — best price-performance, with Chinese open-weight models allowed where the customer explicitly wants them.

Four properties make this enterprise-grade. The posture is set per tenant; it is default-locked to whatever was agreed with the customer; it can be flipped only by an authorized administrator; and every flip is logged. Most routers in the market optimize for cost and quality alone. Almost none let the customer choose — and verify — their own security posture. That is EON's line: best value money can buy, without ever choosing the customer's security posture for them.

7. The model bench, by job

The Switchboard maintains a bench of models mapped to jobs. The table below reflects the landscape as of June 2026; it changes monthly, and the Switchboard absorbs that change so that EON's products never have to.

Job	Best money can buy	Value pick
Code creation (build)	Opus 4.8 / GPT-5.5	Sonnet 4.6
Agentic / terminal	Codex + GPT-5.5 / Claude Code	Kimi K2.7
GUI / frontend	Opus / GPT-5.5	Qwen 3.7 Max
Images	GPT Image 2 (text & labels)	Gemini 3 Pro Image
Inference — quality	GLM-5.2 / DeepSeek V4-Pro	—
Inference — speed / cost	Gemini Flash / Haiku 4.5	DeepSeek V4 Flash

The point is not any single row, which will age. The point is that the customer's product is insulated from a market that moves every week.

8. The economics of routing

Routing converts cost from a function of traffic volume into a function of task complexity. Without routing, spend grows with every request sent to an expensive model. With routing, only the genuinely hard requests reach the frontier, and the rest run on far cheaper tiers. Reported results across the industry cluster between forty and eighty-five percent lower cost, achieved without a visible drop in answer quality — because the bulk of production traffic never required a frontier model in the first place. The routing overhead required to capture those savings is well under a millisecond when decisions are rule-based. For EON, this means the more a customer uses AI — the more value EON delivers — the better the unit economics become, rather than worse.

9. Why this is defensible: the proprietary-data moat

If models are commoditizing, what is defensible? The answer is everything the model is applied to. The equipment ontology, the six-layer competence record, and the orchestration that turns raw model output into reliable work intelligence do not commoditize. Independent analysis points the same way: as high-quality public training data is exhausted over the coming years, the frontier of progress shifts toward companies' proprietary data for specific use cases. That is precisely the Work-Intelligence thesis. The Switchboard is the mechanism that lets EON take cheap, commoditized, interchangeable models and turn them into compounding, owned intelligence. The strategic posture is simple to state: stop renting intelligence; own it instead.

10. Where Switchboard sits: Conductor, Switchboard, Verdict

EON Switchboard is one layer in a three-part control structure, and it is useful to keep the layers distinct.

- **EON Conductor** orchestrates the workflow — what work runs, and in what order.
- **EON Switchboard** decides which model runs each step — for value, quality, and safety.
- **EON Verdict** is the safety gate — it decides whether an output is allowed to ship.

Conductor decides the work. Switchboard decides the model. Verdict decides if it ships. Together they let a customer set a posture once and have every EON product — Genesis, EON Universal, FieldIQ, AssessIQ, Soft Skills, and Customer Virtual — inherit the same value, quality, and safety guarantees.

11. Conclusion

The market is moving toward commoditized models, physical work, and proprietary data. EON Switchboard sits exactly at that intersection. It lets EON serve customers on value, quality, and

safety at the same time — paying frontier prices only where they earn their keep, holding the quality line where it matters, and never compromising a customer's security on their behalf. Understand the work. Route it right. Own the outcome.

Note on currency: model names, benchmark figures, and pricing referenced in this paper reflect the public landscape as of June 2026 and change frequently. The Switchboard architecture is designed precisely to absorb that change. Public materials reference EON's anchor deployments generically as two of the world's largest energy companies.