



Controlled AI Delegation

The Verkflöde Agent Operating Model (VAOM)

Delegation-First · Architecture-Backed · Compliance-Credible

Version 3.5 · 2026

A practical framework for structuring AI decision authority in enterprise workflows.
Confidence gates, delegation boundaries, and audit-ready accountability.

| | |
|--|----|
| Executive Summary | 2 |
| 1. The Delegation Gap | 3 |
| Governance vs Operational Delegation | 4 |
| What Happens Without Delegation Design | 4 |
| 2. What VAOM Is (and What It Is Not) | 4 |
| 3. Core Design Principles | 5 |
| Why This Is Not Workflow Automation | 6 |
| 4. The VAOM Operating Structure | 6 |
| The Confidence Gate: VAOM's Critical Control | 7 |
| 5. Delegation Discovery & Design | 7 |
| 5.1 Decision Inventory | 8 |
| The Problem of Hidden Decisions | 9 |
| 5.2 Authority Decomposition | 10 |
| 5.3 Delegation Pattern Selection | 13 |
| 5.4 Delegation Readiness Assessment | 16 |
| From Discovery to Matrix | 18 |
| 6. Delegation Authority Matrix | 18 |
| 7. Worked Example: Invoice €3,200 | 19 |
| 7A. Worked Example: Customer Complaint Escalation | 20 |
| 7B. Worked Example: AML Transaction Triage | 21 |
| 7C. Worked Example: HR Policy Violation Assessment | 23 |
| 8. How the Composite Confidence Score Works | 25 |
| Implementation Challenges and Calibration Risks | 26 |
| Calibration Anti-Patterns | 27 |
| 9. Managing Foundation Model Drift | 28 |
| 10. Regulatory Alignment | 29 |
| Complementary Frameworks | 30 |
| 11. Implementation Roadmap (90 Days) | 30 |
| Conclusion | 31 |
| Using VAOM | 32 |
| Version History | 33 |

Executive Summary

Enterprise AI is shifting from experimentation to operational delegation. Organizations are no longer testing models in isolation; they are embedding AI agents into workflows that influence financial approvals, contract decisions, HR processes, compliance monitoring, and customer operations. This shift introduces a structural challenge: the management of probabilistic decision-making inside systems historically designed for deterministic automation.

Most enterprises possess governance frameworks covering AI usage, model risk classification, and compliance oversight. But these frameworks define what is allowed in principle. They rarely define how AI systems exercise decision authority within live workflows. The absence of explicit delegation design creates a gap between policy intent and system behavior.

VAOM sits between governance intent and operational execution, translating policy into structured delegation boundaries, confidence thresholds, escalation logic, and audit traceability.

The Verkflöde Agent Operating Model (VAOM) addresses this structural gap. It does not build AI systems. It defines how AI systems are allowed to act within enterprise workflows. VAOM is implementation-agnostic: it works across AI vendors, orchestration platforms, and enterprise systems to establish controlled delegation with measurable outcomes and audit-ready accountability.

VAOM applies wherever AI participates in decisions with operational, financial, or regulatory consequences: invoice processing, contract review, compliance monitoring, HR approvals, customer operations. It is not designed for low-stakes use cases such as internal chatbots or productivity assistants where delegation authority is not a concern.

Critically, VAOM describes what controls must exist, not which products implement them. The seven layers define functional requirements: any combination of orchestration platforms, AI vendors, enterprise systems, and observability tools can satisfy them. This architectural neutrality keeps the framework durable as the underlying technology landscape evolves.

This paper outlines the Delegation Gap, the principles of Delegation-First design, the VAOM layered control structure, a structured method for delegation discovery and design (including six named delegation patterns and a five-dimension authority decomposition framework), a practical worked example, the mechanics and implementation challenges of composite confidence scoring, regulatory alignment, and a 90-day implementation roadmap.

1. The Delegation Gap

Traditional automation executes predefined logic branches. AI agents, by contrast, interpret context and generate probabilistic outputs. When probabilistic systems participate in decision processes, authority must be deliberately structured. Monitoring alone is insufficient; delegation boundaries must be defined before execution occurs.

The Delegation Gap emerges when AI is introduced without redefining authority, escalation, and accountability structures. In practice this manifests as agents that can technically perform actions but have no formal boundaries on when, how, or whether those actions are permitted, or as governance policies that approve AI usage in principle while leaving operational teams without implementable controls.

Governance vs Operational Delegation

Governance defines what is allowed in principle. Delegation defines how AI systems are permitted to act in practice. The distinction is critical and frequently overlooked:

| Governance | Operational Delegation |
|--------------------------------|---|
| Articulates acceptable use | Structures executable authority |
| Classifies risk | Maps decisions to autonomy levels |
| Defines compliance obligations | Embeds obligations into workflow controls |
| Produces documentation | Produces operational architecture |

VAOM does not replace governance. It operationalizes it.

What Happens Without Delegation Design

Consider a scenario that illustrates the Delegation Gap in practice. An AI agent deployed in accounts payable processes vendor invoices. It evaluates payment patterns, matches purchase orders, and detects anomalies. The model is well-trained and performs accurately on standard transactions.

One Friday afternoon, the agent encounters an invoice that references modified payment terms: a vendor has shifted from net-30 to net-15 with a 2% early-payment discount. The agent assesses high confidence: the invoice is legitimate, the vendor is known, the amount is within historical range. It auto-approves the modified terms and posts the payment.

The problem is not that the agent was wrong about the invoice. The problem is that modified payment terms carry contractual and cash-flow implications that fall outside the agent's delegated authority. No one defined that boundary. No confidence gate distinguished between "process this standard invoice" and "accept changed contractual terms." No escalation path existed for this category of decision.

When the CFO discovers the change two weeks later, there is no audit trail explaining why the agent acted, no record of which policy permitted it, and no evidence that a human was ever in the loop. The regulatory inquiry that follows finds the same gap: governance policy approved AI usage in accounts payable, but no operational control defined what the agent was actually allowed to decide.

This is not a model failure. It is a delegation failure. The agent did exactly what it was capable of doing. No one specified what it was allowed to do.

2. What VAOM Is (and What It Is Not)

VAOM is a decision governance framework. It defines what AI systems are allowed to decide within enterprise workflows, under what conditions, with what evidence, and subject to whose authority.

The term "Operating Model" reflects that VAOM governs how decisions operate: the routing logic, the authority boundaries, the confidence thresholds, the escalation paths, and the audit requirements that determine whether a decision is automated, reviewed, or prohibited. It describes the operating logic of delegation, not the operating structure of an organization.

VAOM does not prescribe team structure, reporting lines, or organizational design. It does not define a development lifecycle, a tooling stack, or a deployment methodology. It does not compete with DevOps, MLOps, or platform engineering practices. These are all necessary for building and running AI systems. VAOM addresses a different question: once the system is built and running, what is it allowed to do?

The closest analogy is credit risk decisioning in financial services. Banks have operated for decades with structured decision authority: risk bands that determine which loans can be auto-approved, which require underwriter review, and which must be escalated to credit committee. These systems separate statistical risk scores from institutional authority (a loan officer's approval limit), log every decision with full rationale, and recalibrate thresholds against portfolio performance. VAOM brings this same discipline to AI delegation across all enterprise workflows, not just lending.

This distinction matters because the most common failure mode in enterprise AI governance is not a lack of frameworks. It is frameworks that describe what is allowed in principle while leaving operational teams without implementable controls. VAOM exists to close that gap: not to replace governance, but to make it executable.

3. Core Design Principles

Six principles underpin every VAOM implementation. They ensure that AI participation in workflows enhances performance without eroding governance boundaries.

Controlled Delegation over Blind Automation. Automate only when both statistical confidence and organizational policy permit. Otherwise, route to humans.

Confidence Informs Authority, It Does Not Define It. High confidence does not automatically imply execution rights. Organizational authority (value thresholds, decision categories, regulatory constraints) provides an independent gate.

Architecture Before Acceleration. Define delegation boundaries, controls, and evidence requirements before deploying agent capabilities. Rushing to automation without structure creates ungovernable systems.

Compliance as Structural Outcome. Regulatory requirements (GDPR, DORA, NIS2) are not bolted on after the fact. They are embedded into the operating model so compliance evidence is a byproduct of normal operation.

Human Accountability Preserved. Humans remain accountable for outcomes. The model defines where humans must review, approve, override, or teach, and ensures every human action is logged with the same immutability as agent actions.

Versioned Learning Under Change Control. Every knowledge update, threshold change, or behavior modification is versioned, tested, approved, and auditable. This prevents uncontrolled drift.

Why This Is Not Workflow Automation

Traditional workflow automation, including RPA and rules-based orchestration, executes predefined logic where authority is embedded in the rules themselves. If the rule triggers, the action was already authorized by the person who wrote it. There is no ambiguity about what the system is allowed to do, because it can only do what it was explicitly programmed to do.

AI agents are fundamentally different. They interpret context, weigh evidence, and produce probabilistic outputs that may vary across identical inputs. This means the system exercises something closer to delegated judgment, not scripted execution. The question is no longer "did the rule fire correctly?" but "was the system permitted to make this type of decision at all?"

VAOM exists because that second question has no equivalent in deterministic automation. It requires a control layer that workflow tools were never designed to provide: one that separates statistical confidence from organizational authority and makes delegation boundaries explicit, auditable, and adjustable.

4. The VAOM Operating Structure

VAOM formalizes delegation into seven interconnected control layers. These layers are implementation-agnostic and can be applied across different AI vendors, orchestration platforms, and enterprise systems. Two cross-cutting concerns, Governance and Human Oversight, interact with every layer.

| Layer | Name | Purpose |
|-------|---------------------------------------|--|
| L1 | Trigger & Intake | Captures business events, classifies documents, extracts structured data with PII detection and source validation. |
| L2 | Orchestration | Manages workflow state, routing, retries, escalation paths, and SLA monitoring. The control plane that never makes decisions itself. |
| L3 | Decision & Confidence Gate | Assembles context with provenance, generates a draft decision (never a final one), then evaluates composite confidence against threshold policy to route to auto-execute, review, or escalation. |
| L4 | Controlled Execution | Executes approved actions via secure, idempotent adapters to ERP/CRM/HRIS. Rollback and compensation patterns defined for every write operation. |
| L5 | Knowledge & Learning | Separates read-only business data from tenant-owned agent knowledge. Captures feedback, validates updates against regression suites, and promotes through controlled release. |
| L6 | Governance & Control | Enforces RBAC/ABAC, purpose limitation, data minimization, policy versioning. Creates immutable audit logs, decision traces, metrics, and incident response hooks. |
| L7 | Human Oversight | Review UIs, approval workflows, teaching interfaces, override logging, four-eyes options, and approval SLAs. Humans retain accountability for outcomes. |

The Confidence Gate: VAOM's Critical Control

The Confidence Gate (Layer 3) is the most important and distinctive component in the VAOM architecture. It evaluates a composite confidence score, combining model certainty, rule match strength, data completeness, and anomaly signals, against a threshold policy to route decisions into one of three authority bands:

Band A, Auto-execute: Both statistical confidence and organizational policy permit autonomous action. The agent proceeds without human involvement.

Band B, Human review: Either confidence falls below threshold or the decision category requires oversight. A human reviews and approves before execution.

Band C, Escalation: Novel scenario, high risk, or explicit no-automation zone. Escalated to senior decision-maker or specialist.

Non-delegable: Certain decision types (e.g., contractual modifications, regulatory filings) are never delegated regardless of confidence level. The agent may draft or propose but never execute.

Crucially, the Confidence Gate operates on two independent dimensions: statistical confidence (how certain the system is) and organizational authority (whether the decision category permits automation at all). High confidence alone never grants execution rights. Thresholds are calibrated against historical outcomes and periodically recalibrated through a controlled change process.

5. Delegation Discovery & Design

The seven-layer architecture defines *how* delegation is controlled. The Delegation Authority Matrix captures *what* is delegated. But between architecture and matrix sits a question most frameworks leave unanswered: how does an organization discover which decisions exist, determine which are candidates for delegation, and design the authority boundaries that make delegation safe?

This gap is not theoretical. Organizations adopting AI agents routinely skip delegation design entirely. They select a workflow, connect an agent, and discover the boundaries only when something goes wrong: a modified contract term auto-approved, an escalation path that was never defined, a decision category no one realized the agent could reach. The Delegation Gap described in Section 1 is not just a structural problem. It is a discovery problem. You cannot govern what you have not identified.

VAOM addresses this through a structured Delegation Discovery & Design method, a repeatable process that organizations use to identify decision points, evaluate delegation fitness, design authority boundaries, and produce the Delegation Authority Matrix before any agent capability is deployed.

The method has four stages: Decision Inventory, Authority Decomposition, Delegation Pattern Selection, and Delegation Readiness Assessment. Each stage produces artifacts that feed the next. The final output is a completed Delegation Authority Matrix (Section 6) ready for implementation.

5.1 Decision Inventory

Every workflow contains decisions. Most of them are invisible. An invoice approval workflow does not contain a single decision ("approve or reject"). It contains dozens: classify the document type, extract structured fields, match to a purchase order, validate the vendor, assess whether terms have changed, calculate confidence, determine routing, select the escalation path if confidence is low, and decide what evidence to log. Each of these is a decision point where an AI agent could participate, and each carries different risk, authority, and accountability characteristics.

A Decision Inventory maps every decision point in a target workflow. Not just the primary business decision, but the full decision tree that surrounds it, including the preparatory decisions (data retrieval, context assembly, classification), the routing decisions (which path, which reviewer, which escalation tier), and the meta-decisions (what to log, when to flag, how to handle ambiguity).

How to conduct a Decision Inventory:

Walk the workflow end-to-end with the people who currently execute it. For each step, ask three questions:

What judgment is being applied here? If the answer is "none, it is just a lookup," probe further. A lookup against which data source, selected how, validated by what criteria? Even retrieval involves judgment when AI is involved.

What could go wrong at this step? Each failure mode reveals a decision that is currently being made, often implicitly, about how to handle that failure.

Who is currently accountable for the outcome of this step? If the answer is unclear, the decision point has no owner. That is a delegation design problem regardless of whether AI is involved.

Record each decision point with four attributes:

| Attribute | Description |
|-----------------------------|---|
| Decision ID | Unique identifier within the workflow (e.g., INV-007) |
| Decision description | What judgment is being applied, in plain language |
| Current owner | The role (not person) that currently holds authority |
| Consequence category | Financial, contractual, regulatory, operational, reputational, or informational |

A typical enterprise workflow contains 15 to 40 identifiable decision points. Most organizations, when asked how many decisions their AI agent makes, will name two or three. The inventory reveals the rest.

You cannot delegate authority you have not mapped. The Decision Inventory makes the invisible visible.

The Problem of Hidden Decisions

A Decision Inventory conducted through structured workshops will typically surface 70 to 80 percent of the decision points in a workflow. The remaining 20 to 30 percent are hidden: decisions embedded in tacit knowledge, informal judgment, and undocumented workarounds that the people performing them may not recognize as decisions at all.

These hidden decisions are disproportionately important. They tend to cluster around edge cases, exception handling, and ambiguous situations, precisely the areas where AI delegation is most likely to produce unexpected outcomes. An AP clerk who flags "weird invoices" based on a pattern she cannot articulate is making a decision. If that decision is not inventoried, no delegation boundary will be set for it, and the agent will either replicate her judgment without authorization or ignore it entirely.

Three techniques help surface hidden decisions:

Exception mining. Review the last 6 to 12 months of escalations, overrides, corrections, and rejected transactions. Each exception reveals a decision point that the standard process description omits. If an invoice was manually corrected after auto-processing, a decision was made that the inventory may not have captured.

Shadow observation. Sit with the people who execute the workflow and watch, not for the documented steps, but for the moments where they pause, check something, consult a colleague, or make a judgment call that the process map does not describe. These pauses are hidden decisions.

Adversarial scenario testing. Present the workflow owners with edge cases and ask what would happen: a vendor submitting an invoice in a new currency, a PO that was partially fulfilled, a duplicate that looks legitimate. Each scenario that produces the answer "it depends" or "we would check" reveals a decision point.

No inventory will be complete. The goal is not perfection but coverage sufficient to set delegation boundaries for the decisions that carry material risk. Hidden decisions that are surfaced later, during pilot or production, are captured through the learning pipeline (Layer 5) and fed back into the inventory through the change control process.

5.2 Authority Decomposition

Not every decision in the inventory is a candidate for delegation. Some are trivial and already automated through deterministic rules. Others carry consequences that no organization should delegate to a probabilistic system. Most sit somewhere between.

Authority Decomposition evaluates each decision point against five dimensions that determine its delegation fitness. These dimensions are independent. A decision may score well on four and fail on one, and that single failure may be sufficient to make it non-delegable.

Dimension 1: Reversibility

Can the action be undone? A draft recommendation can be revised. A posted payment is harder to reverse. A regulatory filing, once submitted, may be irreversible. Reversibility determines the cost of being wrong and directly influences whether auto-execution is appropriate.

| Level | Description |
|---------------------------------|---|
| Fully reversible | Action can be undone at negligible cost within a reasonable window (e.g., draft saved but not sent). |
| Reversible with friction | Action can be undone but requires effort, time, or coordination (e.g., ERP posting reversed within 24 hours). |
| Partially reversible | Some consequences can be mitigated but not fully unwound (e.g., payment sent, refund possible but relationship affected). |
| Irreversible | Action cannot be undone (e.g., regulatory filing submitted, contract executed, data deleted). |

Dimension 2: Consequence Scope

Who is affected by this decision, and how broadly? A decision that affects a single internal record has different governance requirements than one that affects a customer, a counterparty, or a regulatory relationship.

| Level | Description |
|-----------------------------|--|
| Internal-operational | Affects internal workflow only (e.g., task routing, queue prioritization). |
| Internal-financial | Affects internal financial position (e.g., budget allocation, payment authorization). |
| External-relational | Affects an external party's experience or relationship (e.g., customer communication, vendor terms). |
| External-regulatory | Affects regulatory standing or creates compliance obligations (e.g., filing, disclosure, reporting). |

Dimension 3: Regulatory Exposure

Does this decision fall within the scope of a specific regulation, and does that regulation impose requirements on how the decision is made? GDPR, DORA, the EU AI Act, NIS2, and sector-specific regulations each create constraints on delegation. A decision that triggers Article 22 of GDPR (automated individual decision-making) has fundamentally different delegation requirements than one that does not.

| Level | Description |
|--|---|
| No direct regulatory exposure | Decision is not individually regulated. |
| General regulatory context | Decision occurs within a regulated process but is not itself a regulated decision. |
| Regulated decision | Specific regulations govern how this decision must be made, documented, or explained. |
| Prohibited from full automation | Regulation requires meaningful human involvement in this specific decision type. |

Dimension 4: Confidence Measurability

Can we actually measure how confident the system is about this decision? Some decisions have clear, quantifiable confidence signals: a document classification score, a match percentage, a statistical certainty. Others involve judgment that resists quantification: whether a contract clause is "materially different," whether a customer communication is "appropriate," whether a risk is "acceptable." If confidence cannot be meaningfully measured, the Confidence Gate in Layer 3 cannot function, and the decision is not a candidate for autonomous execution.

| Level | Description |
|-------------------------------|--|
| Quantifiable | Clear numerical confidence signals exist or can be constructed (e.g., extraction confidence, match score, anomaly probability). |
| Partially quantifiable | Some aspects can be scored but the overall decision involves qualitative judgment (e.g., confidence in data extraction is high, but confidence in interpretation of contract terms is subjective). |
| Judgment-dependent | The decision fundamentally requires qualitative assessment that cannot be reduced to a confidence score without losing its meaning. |

Dimension 5: Accountability Clarity

Is it clear who owns the outcome of this decision? Accountability cannot be delegated to an AI agent. When an agent executes a decision, a human role must remain accountable for the outcome. If the current accountability structure is ambiguous, if no one clearly owns the outcome today, introducing an AI agent will not resolve that ambiguity. It will amplify it.

| Level | Description |
|---|---|
| Clear single owner | One role is unambiguously accountable for this decision outcome. |
| Shared ownership with defined boundaries | Multiple roles share accountability with clear delineation. |
| Ambiguous ownership | Accountability is unclear, disputed, or distributed without clear boundaries. |
| No defined owner | Nobody currently owns this decision outcome explicitly. |

Record the assessment for each decision point:

| Decision ID | Reversibility | Consequence Scope | Regulatory Exposure | Confidence Measurability | Accountability Clarity |
|-------------|--------------------------|----------------------|---------------------|--------------------------|------------------------|
| INV-001 | Fully reversible | Internal-operational | No direct exposure | Quantifiable | Clear single owner |
| INV-007 | Reversible with friction | Internal-financial | General context | Quantifiable | Clear single owner |
| INV-012 | Irreversible | External-regulatory | Regulated decision | Judgment-dependent | Ambiguous ownership |

The decomposition does not produce a single score. It produces a profile. A decision that is fully reversible, internally scoped, unregulated, quantifiable, and clearly owned is a strong delegation candidate. A decision that is irreversible, externally scoped, regulated, judgment-dependent, and ambiguously owned is almost certainly non-delegable. Most decisions fall between these extremes, and the profile guides where they land in the Delegation Authority Matrix.

Authority Decomposition separates "can the AI do this?" from "should the AI be allowed to do this?" The first question is technical. The second is organizational. VAOM answers the second.

5.3 Delegation Pattern Selection

Once the Decision Inventory has mapped the decision landscape and Authority Decomposition has profiled each decision's delegation fitness, the next step is to assign each delegable decision to a

delegation pattern: a named, reusable configuration of agent authority, human involvement, and evidence requirements.

VAOM defines six delegation patterns. Each pattern represents a distinct relationship between agent action and human authority. The patterns are ordered by increasing agent autonomy.

Pattern 1: Prepare & Present

The agent assembles context, retrieves data, and organizes information for human decision-making. The agent makes no decision and takes no action. The human decides and acts.

Use when: The decision is high-consequence, judgment-dependent, or non-delegable, but the preparation work is time-consuming and suitable for automation.

Human role: Decision-maker. *Agent role:* Context assembler. *Evidence:* What the agent retrieved, how it was assembled, what was presented.

Pattern 2: Draft & Approve

The agent assembles context *and* produces a draft decision or recommendation. The human reviews, modifies if necessary, and approves before execution.

Use when: The decision requires human judgment but the agent can reliably produce a reasonable first draft that reduces human effort.

Human role: Approver with modification rights. *Agent role:* Drafter. *Evidence:* Draft produced, human modifications made, approval timestamp, approver identity.

Pattern 3: Triage & Route

The agent classifies incoming items and routes them to the appropriate handler (human or automated) based on defined rules and confidence thresholds. The agent does not resolve the item; it determines who or what should.

Use when: Volume is high, classification criteria are well-defined, and the cost of misrouting is manageable (items reach a human who can re-route).

Human role: Exception handler for misrouted or unclassifiable items. *Agent role:* Classifier and router. *Evidence:* Classification rationale, confidence score, routing decision, exception handling log.

Pattern 4: Execute & Audit

The agent makes the decision and executes the action autonomously. Humans review a sample of completed decisions after the fact through scheduled audits.

Use when: Confidence is high, consequences are bounded and reversible, organizational authority permits autonomous action, and historical data supports reliable threshold calibration.

Human role: Auditor (post-execution review of samples and exceptions). *Agent role:* Decision-maker and executor within defined authority. *Evidence:* Full decision trace, confidence score, authority verification, audit sample selection criteria, audit findings.

Pattern 5: Monitor & Intervene

The agent operates continuously, monitoring, detecting, and responding to events within defined parameters. Humans are alerted when thresholds are crossed or anomalies detected, and retain the ability to intervene at any point.

Use when: Continuous operation is required (e.g., compliance monitoring, anomaly detection, SLA tracking), the agent's detection capabilities exceed human capacity at scale, and intervention paths are well-defined.

Human role: Interventionist (responds to alerts, overrides when needed). *Agent role:* Continuous monitor with bounded response authority. *Evidence:* Monitoring logs, alert triggers, intervention records, false-positive rates, response times.

Pattern 6: Coordinate & Escalate (Multi-Agent)

Multiple agents collaborate on a workflow, with each agent operating within its own delegation boundaries. A coordinating agent or orchestration layer manages handoffs. Escalation to humans occurs when any agent in the chain encounters a decision outside its authority, when the aggregate confidence across the chain falls below threshold, or when the coordination itself produces an unexpected state.

Use when: The workflow spans multiple domains or systems, specialized agents handle different aspects, and the organization has mature governance over individual agent authorities.

Human role: Escalation authority for cross-boundary decisions and coordination failures. *Agent role:* Specialist within defined scope, with explicit handoff protocols. *Evidence:* Per-agent decision traces, handoff logs, coordination state, escalation triggers, end-to-end audit trail across the full agent chain.

Multi-agent coordination introduces three failure modes that single-agent patterns do not face:

Authority conflicts. Two agents may claim jurisdiction over the same decision, or an agent may receive a handoff that falls outside its delegation boundary but inside no other agent's boundary either. The coordination layer must define a conflict resolution protocol: which agent takes precedence, under what conditions, and what happens when no agent has authority. Unresolved authority conflicts must escalate to a human, not be silently resolved by the coordination layer.

Cascading confidence erosion. When agents operate in sequence, each agent's output becomes the next agent's input. Small errors or low-confidence outputs early in the chain can compound: an extraction agent that misreads a field at 0.88 confidence feeds a validation agent that passes it at 0.91 confidence, which feeds a decision agent that auto-approves at 0.93 confidence. Each individual score looks acceptable; the aggregate reliability may not be. Multi-agent chains should track cumulative confidence across the full chain, not just per-agent scores. If the product of individual confidence scores falls below a chain-level threshold, the entire decision routes to human review regardless of individual scores.

Coordination state loss. When handoffs between agents lose context, such as a field that was flagged as uncertain by one agent but not passed through the handoff protocol, downstream agents make decisions on incomplete information without knowing it is incomplete. Handoff protocols must define which state is transferred, which is dropped, and what metadata (including uncertainty markers and provenance) must survive the transition.

Each decision point from the inventory is assigned to a pattern based on its Authority Decomposition profile:

| Decomposition Profile | Typical Pattern |
|--|-----------------------|
| Irreversible + regulated + judgment-dependent | Prepare & Present |
| High-consequence + quantifiable + clear owner | Draft & Approve |
| High-volume + well-defined criteria + manageable misroute cost | Triage & Route |
| Bounded + reversible + high confidence + clear authority | Execute & Audit |
| Continuous + detection-oriented + defined intervention paths | Monitor & Intervene |
| Multi-domain + multiple specialized agents + mature governance | Coordinate & Escalate |

Delegation patterns give organizations a shared vocabulary for delegation design. "We use Execute & Audit for standard invoice approvals and Draft & Approve for contract-adjacent decisions" communicates more in one sentence than a page of policy documentation.

5.4 Delegation Readiness Assessment

Before a delegation design moves to implementation, each decision point assigned to a delegation pattern must pass a readiness assessment. This is not a maturity model. It is a go/no-go checklist for a specific decision in a specific workflow. A decision that is not ready for delegation is not ready, regardless of organizational AI maturity or agent capability.

The assessment evaluates five readiness conditions. All five must be met for the delegation to proceed. A failure on any single condition means the delegation design is incomplete and must be resolved before implementation.

Condition 1: Measurable confidence signals exist. The Confidence Gate (Layer 3) requires quantifiable inputs. For this specific decision type, can the system produce a composite confidence score from model certainty, rule match strength, data completeness, and anomaly signals? If confidence cannot be measured, Band A (auto-execute) routing is not possible and the delegation pattern must account for this. *Evidence required:* Documentation of which confidence signals are available, how they combine into a composite score, and what historical data supports the threshold calibration.

Condition 2: Authority boundaries are explicit. The Delegation Authority Matrix entry for this decision must define: which authority band applies at each confidence level, what value or risk thresholds modify the routing, and which decision categories are non-delegable regardless of confidence. Boundaries cannot be implicit or assumed. *Evidence required:* Completed matrix entry, signed off by the role that holds authority for the decision type.

Condition 3: Escalation paths are defined and tested. When the agent encounters a decision outside its authority (confidence too low, decision category requires human review, novel scenario detected), what happens? The escalation path must be defined, the receiving role must be identified, and the path must be tested to confirm that escalation actually reaches a human with the authority and context to act. *Evidence required:* Documented escalation path, identified receiving role, test results showing escalation completes within defined SLA.

Condition 4: Audit evidence is producible. If a regulator, auditor, or senior stakeholder asks "why did the system make this decision?", can the organization answer? The answer requires: the decision ID, the inputs considered, the confidence score, the policy that permitted the action, the authority band that applied, the execution timestamp, and the identity of the accountable human. If any of these cannot be produced, the evidence chain is broken. *Evidence required:* Sample audit record for a test decision, demonstrating all required evidence fields are populated.

Condition 5: A named human is accountable. A specific role (not a team, not a committee, not "management") is identified as accountable for outcomes produced by this delegation. That role has the authority to override, suspend, or modify the delegation. The role is documented, and the person in it knows they hold this accountability. *Evidence required:* Named accountable role in the delegation design, confirmation that the role holder has accepted accountability, override mechanism documented.

In practice, satisfying Condition 5 is often the hardest part of the readiness assessment. Authority in enterprises is frequently contested, shared across functions, or deliberately left ambiguous. Invoice approval rules may sit between Finance, Procurement, and Legal. Vendor exceptions may be owned regionally in some contexts and globally in others. Nobody may want to own AI-delegated decisions because accountability for automated outcomes feels riskier than accountability for manual ones.

This is not a problem VAOM can solve by design. It is a problem VAOM makes visible. If no one is willing to be named as accountable for a delegated decision, that decision is not ready for delegation, regardless of how technically capable the agent is. The readiness assessment forces the conversation that the organization may have been avoiding. Three approaches help navigate contested ownership: escalate the accountability question to the executive sponsor before the design phase begins (not during it), frame accountability as override authority rather than blame assignment (the accountable role is the one that can suspend the delegation, not the one that gets fired if it fails), and start with decision types where ownership is already clear, building trust before attempting decisions where authority is disputed.

The Delegation Readiness Assessment produces one of three outcomes:

Ready. All five conditions met. Proceed to implementation (Phase 2 onward in the 90-day roadmap).

Conditionally ready. One or two conditions partially met with a clear remediation path and timeline. Proceed to implementation with the remediation built into the plan.

Not ready. One or more conditions fundamentally unmet. Do not proceed. Resolve the gap before revisiting.

Readiness is decision-specific, not organization-wide. An organization may be ready to delegate invoice classification (Pattern 3: Triage & Route) while being entirely unready to delegate contract term assessment (which may require Pattern 1: Prepare & Present until confidence signals improve). This specificity prevents both over-caution and over-delegation.

From Discovery to Matrix

The four stages of Delegation Discovery & Design feed directly into the Delegation Authority Matrix (Section 6). The inventory identifies decision points. The decomposition profiles their delegation fitness. The patterns define the agent-human relationship. The readiness assessment confirms the design is implementable.

The resulting matrix is not a theoretical artifact. It is a governance decision record: explicit, auditable, and owned by a named authority. It can be reviewed, challenged, revised, and version-controlled. And critically, it can be communicated: to the team implementing the agent, to the compliance function validating the controls, to the regulator examining the evidence, and to the executive accountable for the outcome.

This is what it means to operationalize governance. Not a policy that describes what is allowed in principle, but a design process that produces implementable delegation boundaries, backed by evidence, owned by named humans, and ready for the seven-layer control architecture to enforce.

6. Delegation Authority Matrix

The matrix below is the output of the Delegation Discovery & Design process described in Section 5. Each row represents a decision point identified in the Decision Inventory, profiled through Authority Decomposition, and assigned a delegation pattern.

A Delegation Authority Matrix maps specific decision types to permitted levels of autonomy under defined confidence and risk conditions. This matrix makes delegation explicit, auditable, and adjustable.

Example: Vendor Invoice Approval Workflow

| Decision Type | High Confidence | Medium Confidence | Low / No Confidence |
|-----------------------------|-----------------|-------------------|---------------------|
| Standard invoice ≤ €5k | ✔ Auto-approve | 👁 Human review | ⚠ Escalation |
| Invoice > €5k | 👁 Human review | 👁 Human review | ⚠ Escalation |
| Duplicate detection anomaly | 👁 Human review | ⚠ Escalation | ⚠ Escalation |
| Contractual modification | 🚫 Non-delegable | 🚫 Non-delegable | 🚫 Non-delegable |

Each row maps to a delegation pattern from Section 5.3. Standard invoices below €5k use Execute & Audit (Pattern 4): the agent decides and acts, humans audit a sample. Invoices above €5k use Draft & Approve (Pattern 2): the agent produces a recommendation, a human approves before execution. Duplicate detection anomalies use Triage & Route (Pattern 3): the agent flags and routes, a human investigates. Contractual modifications are non-delegable; the agent may use Prepare & Present (Pattern 1) to assemble context, but a human makes and executes the decision.

This matrix separates statistical confidence from organizational authority. High confidence does not automatically imply execution rights.

7. Worked Example: Invoice €3,200

To illustrate VAOM in practice, we trace a single vendor invoice through the delegation discovery process and all seven control layers. The invoice workflow was selected through the Decision Inventory, which identified 23 decision points across the accounts payable process. The €5k threshold emerged from Authority Decomposition (consequence scope analysis). The auto-approve routing for standard invoices below €5k uses the Execute & Audit pattern (Pattern 4). This demonstrates how each layer contributes specific controls, evidence, and decision logic.

Layer 1, Trigger & Intake. A vendor invoice email arrives. Document Intelligence classifies it as invoice/standard, extracts the amount (€3,200), PO reference, line items, and VAT. PII is redacted. Per-field confidence scores are attached.

Layer 2, Orchestration. The orchestrator creates a task record, enriches it with a PO match from the ERP, and routes to the Decision layer. Retry policy: 2x on ERP timeout with exponential backoff. SLA timer starts.

Layer 3, Decision & Confidence Gate. Context Assembly retrieves vendor history, contract terms, and 3-way match data with full provenance. The Decision Proposal applies rules and model reasoning to produce a draft: approve, PO matched, within contract terms, below €5k threshold. The Confidence Gate evaluates a composite score of 0.94 (high band). Since the invoice is below €5k and confidence is high, Band A routing is selected: auto-approve.

Layer 4, Controlled Execution. The approval is posted to the ERP via an idempotent API call. Transaction ID is logged. Rollback path: reversal available within 24 hours.

Layer 5, Knowledge & Learning. Invoice data remains in the read-only ERP store. Approval rationale is written to the tenant-owned knowledge store (EU-resident, encrypted). A correction on a similar invoice last month had already been distilled into a threshold refinement (0.90 to 0.92), validated against 200 historical invoices, and promoted.

Layer 6, Governance & Control. RBAC verified: the agent holds invoice.approve permission for ≤€5k, purpose-limited to the AP process. An immutable audit log records the decision ID, rationale, confidence score, all data accessed, policy checks passed, and execution timestamp.

Layer 7, Human Oversight. No human review is needed (high confidence, within authority). A 5% sampling audit is scheduled. The override UI remains available should any stakeholder wish to intervene after the fact.

Total time from email arrival to ERP posting: under 90 seconds. Audit evidence available within minutes of a regulatory inquiry.

7A. Worked Example: Customer Complaint Escalation

To demonstrate VAOM beyond financial workflows, we trace a customer complaint through a retail banking operations context. The complaint workflow was selected through the Decision Inventory, which identified 18 decision points. The delegation design uses three patterns across different decision types: Triage & Route (Pattern 3) for initial classification, Draft & Approve (Pattern 2) for response generation, and Prepare & Present (Pattern 1) for regulatory reporting decisions.

Layer 1, Trigger & Intake. A customer submits a complaint via the bank's online portal about unauthorized charges on their account. The intake system captures the complaint text, customer ID, account details, and transaction references. PII is tagged and access-restricted. Sentiment analysis attaches an urgency score.

Layer 2, Orchestration. The orchestrator creates a case record, enriches it with the customer's account history, prior complaints, and product holdings. Regulatory SLA timer starts: the complaint must be acknowledged within 24 hours and resolved or escalated within 15 business days under FCA requirements.

Layer 3, Decision & Confidence Gate. Three decision points pass through the gate in sequence:

Classification decision. The agent classifies the complaint as "unauthorized transaction, potential fraud" with 0.91 confidence. This classification determines which team handles the case and whether fraud investigation is triggered. The Confidence Gate routes this as Band A: the classification criteria are well-defined and the cost of misrouting is manageable (the receiving team can re-route).

Response draft decision. The agent drafts an acknowledgment email and a preliminary assessment. Because customer-facing communications carry reputational and regulatory consequences (consequence scope: external-relational), this decision uses Draft & Approve regardless of confidence. The draft is queued for human review.

Regulatory reporting decision. The agent assesses whether the complaint triggers mandatory regulatory reporting. This is non-delegable: the agent assembles the relevant data and flags the regulatory criteria that may apply, but a compliance officer makes the determination.

Layer 4, Controlled Execution. The classification is recorded in the CRM. The acknowledgment email is sent only after human approval. The case is routed to the fraud investigation queue with full context attached.

Layer 5, Knowledge & Learning. The complaint joins the anonymized complaint corpus used for classification training. A pattern of similar unauthorized-charge complaints over the past month is flagged for the fraud team's attention, not as an automated action but as a data signal surfaced through the monitoring pipeline.

Layer 6, Governance & Control. The agent's CRM access is scoped to complaint-related records only. The regulatory reporting flag is logged with the compliance officer's determination, creating an audit trail for FCA inquiries. All customer data access is purpose-limited and retention-scheduled.

Layer 7, Human Oversight. The complaint handler reviews the draft response, modifies the tone to address the customer's specific frustration, and approves sending. The compliance officer reviews the regulatory reporting assessment and determines that no FCA report is required for this individual complaint but notes the pattern for quarterly review.

This example illustrates how different decision types within a single workflow can use different delegation patterns simultaneously, and how consequence scope (not just confidence) determines the routing.

7B. Worked Example: AML Transaction Triage

Anti-money laundering (AML) transaction monitoring operates at a scale and regulatory intensity that makes it an instructive VAOM application. A mid-sized European bank processes approximately 2 million transactions daily. The current AML system generates around 800 alerts per day, of which historical analysis shows roughly 95 percent are false positives. Human investigators spend the majority of their time dismissing alerts that should never have reached them.

The Decision Inventory identified 12 decision points in the alert triage workflow. The delegation design uses Monitor & Intervene (Pattern 5) for continuous transaction screening, Triage & Route (Pattern 3) for alert prioritization, and Prepare & Present (Pattern 1) for suspicious activity report (SAR) preparation.

Layer 1, Trigger & Intake. The transaction monitoring system generates an alert: a series of structured cash deposits just below the reporting threshold from an account that has historically shown only salary and utility activity. The alert includes transaction details, account profile, historical patterns, and the rule that triggered the flag.

Layer 2, Orchestration. The orchestrator assigns the alert a priority score based on regulatory risk (structuring patterns are high-priority under the Fourth Anti-Money Laundering Directive). SLA timer starts: high-priority alerts must be triaged within 4 hours.

Layer 3, Decision & Confidence Gate. Two decision points pass through the gate:

Triage decision. The agent evaluates the alert against known typologies, account history, and contextual signals. For this alert, rule match strength is high (the deposit pattern matches structuring typology), model certainty is 0.87 (the account holder has no prior suspicious activity, which introduces ambiguity), and data completeness is strong (all transaction records are available). Composite confidence: 0.83, routing to Band B (human review). The agent cannot dismiss this alert, but it can prioritize it and assemble the investigation package.

Dismissal decision for low-risk alerts. Separately, the agent evaluates a batch of 200 alerts that match known false-positive patterns: recurring transfers between the customer's own accounts, salary payments matching employer records, and utility payments within historical range. For these, all four confidence dimensions score above 0.95. Band A routing permits auto-dismissal with full logging. This is where the operational value concentrates: reducing the 800 daily alerts to approximately 40 to 60 that require human investigation.

Layer 4, Controlled Execution. Auto-dismissed alerts are closed in the case management system with a coded rationale and full decision trace. The structuring alert is escalated to the investigation queue with the agent's assembled context package.

Layer 5, Knowledge & Learning. Investigator outcomes feed back into the system: confirmed false positives refine the auto-dismissal criteria; confirmed suspicious activity refines the typology models. Every refinement follows the change control process. A recent regulatory guidance update on virtual asset service providers was incorporated as a new typology rule, validated against 6 months of historical alerts, and promoted after compliance sign-off.

Layer 6, Governance & Control. AML regulatory requirements impose specific constraints: no alert may be auto-dismissed if it involves a politically exposed person (PEP), a sanctioned jurisdiction, or an amount exceeding the regulatory reporting threshold. These are hard-coded as non-delegable overrides in the Confidence Gate policy, regardless of composite score. The audit log captures every dismissal rationale and every escalation, producing the evidence trail required under the Fourth Anti-Money Laundering Directive.

Layer 7, Human Oversight. The AML investigator receives the structuring alert with the full context package: transaction timeline, account profile, typology match analysis, and the agent's confidence breakdown. The investigator determines that the deposits correlate with the account holder's recently registered cash-intensive small business, documented in a KYC update from three months prior. The alert is closed as a false positive, and the investigator's rationale is logged. The auto-dismissal sample audit (10 percent of auto-dismissed alerts reviewed weekly) confirms no missed true positives this period.

This example demonstrates VAOM in a heavily regulated, high-volume context where the primary value is not faster decisions but better-targeted human attention: reducing noise so investigators focus on the alerts that matter.

7C. Worked Example: HR Policy Violation Assessment

HR workflows involve decisions that affect individuals' careers and livelihoods, making them a demanding test for delegation design. A multinational employer uses an AI agent to assist with policy violation assessments, following a complaint about a manager's conduct during a team meeting.

The Decision Inventory identified 14 decision points in the policy violation workflow. The critical insight from Authority Decomposition was that nearly every decision in this workflow scores high on consequence scope (external-relational: affects a person's employment relationship) and accountability clarity is often ambiguous (HR, Legal, the line manager's manager, and sometimes a works council all have a role). The delegation design uses Prepare & Present (Pattern 1) for the majority of decisions, with Triage & Route (Pattern 3) only for the initial intake classification.

Layer 1, Trigger & Intake. An employee submits a complaint through the HR case management system, reporting that a manager made demeaning remarks during a team meeting. The intake system captures the complaint text, identifies the parties involved, and flags the applicable policies (dignity at work, anti-harassment). PII protections are elevated: access is restricted to the HR case handler and the designated investigator.

Layer 2, Orchestration. The orchestrator creates a case record with restricted visibility. It enriches the case with prior complaints involving either party (if any exist) and the applicable local employment law requirements. In the EU, works council notification requirements are checked. SLA timer starts per internal policy.

Layer 3, Decision & Confidence Gate. Three decision points are evaluated:

Classification decision. The agent classifies the complaint as "conduct, dignity at work, severity: moderate" with 0.86 confidence. Because misclassification could result in an investigation being scoped too narrowly or too broadly, and because the consequence scope is external-relational, this routes to Band B: a human HR case handler reviews the classification before the investigation is scoped.

Evidence sufficiency decision. The agent assesses whether the complaint contains enough information to proceed to investigation or whether further information is needed. This is Prepare & Present: the agent drafts a summary of what is known and what gaps exist, but the HR case handler decides whether to request additional information.

Outcome recommendation. This is non-delegable. The agent does not recommend disciplinary outcomes. It assembles the evidence package: complaint details, witness statements (if gathered by the investigator), applicable policy provisions, precedent from prior similar cases (anonymized), and local legal requirements. The investigator and HR decision-maker determine the outcome.

Layer 4, Controlled Execution. The only automated execution is the case record creation and the notification to the designated HR case handler. All subsequent actions (investigation initiation, witness outreach, outcome determination) are human-executed, with the agent providing preparation and context assembly.

Layer 5, Knowledge & Learning. Case outcomes are stored in an anonymized, access-restricted knowledge base used to improve classification accuracy and to surface patterns (e.g., repeated complaints about the same team or location). Learning updates to the classification model follow enhanced change control: HR Legal must approve any changes to how complaints are categorized, given the employment law implications.

Layer 6, Governance & Control. Access controls are unusually strict: the case is visible only to the assigned HR case handler, the designated investigator, and HR Legal. The agent's data retrieval is purpose-limited to the specific case. GDPR data subject rights apply: both the complainant and the respondent have rights regarding their personal data in the case file. The audit log records every data access, every document generated, and every human decision, but the log itself is access-restricted.

Layer 7, Human Oversight. Humans are in the loop at every substantive decision. The agent's value is in preparation and pattern recognition, not decision-making. The HR case handler reviews the classification, scopes the investigation, and determines the process. The investigator conducts interviews and gathers evidence. The decision-maker (senior HR with Legal input) determines the outcome. The agent never recommends discipline, never contacts witnesses, and never communicates outcomes.

This example illustrates a workflow where VAOM's primary contribution is establishing clear non-delegable boundaries. The value of the framework here is not automation but clarity: defining precisely what the agent may and may not do in a context where the consequences of over-delegation are severe.

8. How the Composite Confidence Score Works

The Confidence Gate (Section 4, Layer 3) references a composite confidence score but does not specify how it is constructed. This section defines the scoring model.

The composite score combines four dimensions, each contributing a distinct signal:

Model certainty: The AI model's own probability or logit-based confidence in its output. This is the most intuitive dimension but also the least reliable in isolation, as model confidence can be poorly calibrated.

Rule match strength: The degree to which deterministic business rules confirm or contradict the model's output. A high rule match (e.g., PO number matches exactly, amount within contract terms) reinforces confidence. A rule mismatch (e.g., vendor not in approved list) reduces it regardless of model certainty.

Data completeness: Whether all required input fields are present, validated, and within expected ranges. Missing or malformed data reduces confidence even when the model produces a high-certainty output from incomplete inputs.

Anomaly signals: Deviation from historical patterns. Unusual amounts, unfamiliar vendors, atypical timing, or format anomalies act as negative modifiers. These signals can demote an otherwise high-confidence decision to human review.

The composite score is evaluated using a configurable policy that defines how dimensions combine. In the simplest implementation, this is a weighted average with a hard floor: if any single dimension falls below a minimum threshold, the composite score is capped at the review band regardless of other signals. More sophisticated implementations may use Boolean AND logic for critical dimensions (e.g., data completeness must always pass) combined with weighted scoring for probabilistic dimensions.

The key design principle is that the scoring policy is explicit, versioned, and auditable. The organization defines the weights, the floors, and the combination logic. The gate does not rely on opaque model internals. When a regulator asks "why did the system auto-approve this?", the answer is traceable to specific scores on specific dimensions against a specific policy version, not to a black-box confidence number.

Implementation Challenges and Calibration Risks

The composite confidence model is architecturally sound but non-trivial to implement well. Three challenges deserve honest acknowledgment:

Model confidence is often poorly calibrated. Foundation models frequently produce high-confidence outputs that are incorrect and low-confidence outputs that are correct. A model that reports 0.95 certainty on an invoice classification may be wrong 15 percent of the time at that threshold, not 5 percent. This is why VAOM treats model certainty as one of four dimensions rather than the sole input, and why the hard floor mechanism exists: a single unreliable signal cannot override the composite. Organizations implementing the Confidence Gate should calibrate model certainty against ground truth data for their specific decision types, not rely on the model's self-reported probabilities.

Anomaly detection generates noise. Anomaly signals are powerful negative modifiers, but real-world anomaly detectors produce false positives. An unusual invoice timing may reflect a vendor's changed billing cycle, not fraud. A new vendor format may simply be a new vendor. If anomaly signals are weighted too heavily, the Confidence Gate routes an excessive proportion of decisions to human review, creating the bottleneck problem: review queues grow, SLAs break, and humans become slower than the process was before the agent was introduced. Threshold calibration must account for the false positive rate of each anomaly signal, and organizations should expect to recalibrate anomaly weights multiple times during the pilot phase.

Rules and model reasoning can conflict. When deterministic business rules contradict the model's probabilistic assessment, the composite score must resolve the conflict. VAOM's design gives rule match strength its own dimension precisely for this reason, but the resolution policy must be explicit. In most implementations, a hard rule failure (vendor not on approved list, amount exceeds contract ceiling) should cap the composite score at Band B or below regardless of model certainty. The alternative, letting high model confidence override a failed rule, defeats the purpose of having rules.

The Confidence Gate is not a plug-and-play component. It requires investment in calibration data, iterative threshold tuning, and ongoing monitoring of score distributions. The framework provides the architecture; the organization must provide the engineering discipline to make it reliable. When implemented well, it produces measurably better outcomes than either pure automation or pure human review. When implemented carelessly, with uncalibrated thresholds and unweighted dimensions, it becomes either a bottleneck that routes everything to humans or a rubber stamp that auto-approves decisions it should not.

Calibration Anti-Patterns

Miscalibrated confidence thresholds are the most common cause of VAOM implementation failure. These anti-patterns represent failures of delegation design and governance discipline, not failures of AI capability. The architecture works correctly, the layers are wired, the audit logs are flowing, but the thresholds are wrong, and the system either produces no value (everything goes to humans) or produces uncontrolled risk (too much is auto-approved). The following anti-patterns describe what bad calibration looks like and what signals indicate recalibration is needed.

Anti-pattern 1: The Review Queue Flood. Symptom: 60 to 80 percent of decisions route to Band B (human review). Review queues grow. SLAs degrade. Humans begin rubber-stamping approvals to clear the backlog, which is worse than not having the system at all because rubber-stamping creates the illusion of oversight without its substance. Root cause: thresholds set too conservatively, often because the implementation team defaulted to "safe" settings without calibrating against historical data. The anomaly dimension may also be over-weighted, treating normal business variation as suspicious. Signal to watch: if the human override rate on Band B decisions exceeds 90 percent (meaning humans approve more than 90 percent of what the agent flagged for review), the threshold is too strict. The agent is not adding judgment; it is adding delay.

Anti-pattern 2: The Confidence Mirage. Symptom: the system auto-approves at high rates and the composite scores look healthy, but post-execution audits reveal error rates significantly above expectations. Root cause: over-reliance on model certainty without adequate weight on rule match and data completeness. The model reports high confidence, but the confidence is poorly calibrated for this specific decision type. The composite score inherits the model's overconfidence because the weighting policy gives model certainty too much influence. Signal to watch: if the composite score distribution clusters tightly above the Band A threshold (e.g., 85 percent of scores fall between 0.92 and 0.98), the scoring lacks discrimination. A healthy distribution shows meaningful spread across the bands, reflecting genuine variation in decision difficulty.

Anti-pattern 3: The Exception Graveyard. Symptom: the Delegation Authority Matrix has accumulated a long list of hardcoded exceptions, special cases, and manual overrides that bypass the Confidence Gate. Each exception was added to handle a specific situation, but collectively they have created a shadow routing system that the Confidence Gate does not govern. Root cause: instead of recalibrating thresholds or adjusting dimension weights when the gate produces incorrect routings, the team adds exceptions. Each exception is individually reasonable, but the aggregate effect is that the Confidence Gate governs a shrinking proportion of decisions while the exception list grows. Signal to

watch: if more than 15 percent of decisions are routed through exception rules rather than through the standard composite scoring path, the scoring policy itself needs revision, not more exceptions.

Anti-pattern 4: The Stale Threshold. Symptom: the system performed well during pilot but accuracy has degraded over 3 to 6 months of production operation. Root cause: the business has changed (new vendors, new pricing structures, seasonal patterns, updated regulations) but the thresholds, dimension weights, and anomaly baselines have not been recalibrated. The Confidence Gate is evaluating current decisions against historical calibration that no longer reflects the operating environment. Signal to watch: a gradual shift in the band distribution over time. If the percentage of Band A decisions is drifting upward or downward without a corresponding change in business volume or complexity, the thresholds are drifting out of alignment with reality.

Anti-pattern 5: The Dimension Collapse. Symptom: the composite score behaves almost identically to a single dimension, usually model certainty. The other three dimensions (rule match, data completeness, anomaly signals) are either not implemented, not weighted meaningfully, or always return high scores. Root cause: the implementation team built the model certainty pipeline first and either deferred the other dimensions or implemented them as placeholders that return constant values. The composite score exists in name but not in practice. Signal to watch: calculate the correlation between the composite score and each individual dimension. If any single dimension has a correlation above 0.95 with the composite, the other dimensions are not contributing meaningful signal and the gate is effectively single-dimensional.

These anti-patterns are not theoretical. They emerge in production systems within the first 3 to 12 months of operation. Organizations implementing the Confidence Gate should monitor for these signals as part of the ongoing governance process described in Layer 6, and should schedule formal recalibration reviews at least quarterly during the first year of operation.

9. Managing Foundation Model Drift

A distinct challenge for AI-enabled workflows is foundation model drift: the risk that an underlying AI model changes, through provider updates, version deprecation, or fine-tuning adjustments, in ways that invalidate previously calibrated confidence thresholds and delegation boundaries.

VAOM addresses this through Layer 5 (Knowledge & Learning) and Layer 6 (Governance & Control) working together. When a foundation model change is detected or announced, the following control sequence applies:

Regression testing. The updated model is evaluated against the same historical decision set used to calibrate the original thresholds. Outputs are compared for consistency, confidence distribution shifts, and decision boundary changes.

Threshold recalibration. If regression testing reveals material drift (confidence distributions shifted, previously auto-approved decisions now falling below threshold, or vice versa), thresholds are recalibrated against the new model's behavior. Recalibration follows the same change control process as any threshold modification.

Shadow period. Before the updated model enters production, it runs in parallel with the existing model for a defined period. Decisions are logged but not executed. Divergences between old and new model outputs are reviewed by the designated human authority.

Model registry. Every model version used in production is recorded in a registry with its calibration data, regression test results, shadow period outcomes, and the human approval that authorized its promotion. This creates an auditable chain from model version to threshold policy to delegation authority.

The principle is straightforward: a model change is a change to the delegation system. It is governed by the same change control discipline as any other modification to the operating model. Provider convenience does not override organizational governance.

10. Regulatory Alignment

VAOM embeds compliance into operational design so that regulatory evidence is a byproduct of normal operation, not a separate workstream. The model maps to three key European regulatory frameworks:

| Regulation | Key Requirements | VAOM Layers Providing Evidence |
|-------------|--|--|
| GDPR | Data minimization, purpose limitation, transparent decision logic, right to explanation | L1 (PII detection, minimization at entry), L3 (purpose-limited retrieval, traceable reasoning), L6 (immutable decision logs, access audit trail) |
| DORA | Operational resilience, change control, third-party risk management, audit evidence generation | L2 (retry/escalation/SLA), L4 (idempotent execution, rollback), L5 (versioned learning, regression tests), L6 (evidence packs, incident runbooks) |
| NIS2 | Security governance, monitoring, incident response, accountability | L6 (RBAC/ABAC, policy versioning, real-time alerting, SLO monitoring), L7 (override logging, escalation paths), L4 (secrets management, zero-trust auth) |

Complementary Frameworks

VAOM does not replace existing enterprise frameworks. It complements them by operationalizing AI-specific delegation controls within structures that organizations already maintain:

EU AI Act: VAOM's confidence gate and human oversight layers map directly to the Act's requirements for high-risk AI systems, particularly human oversight (Article 14), transparency (Article 13), and risk management (Article 9). The Delegation Discovery method (Section 5) provides the structured process for assessing whether a decision falls within the Act's scope and designing appropriate controls.

NIST AI RMF: VAOM implements the Govern, Map, Measure, and Manage functions at the workflow level, providing operational specificity where the NIST framework provides strategic guidance.

ISO 42001: VAOM's governance layer (L6) aligns with the management system requirements, while the learning pipeline (L5) addresses continual improvement.

11. Implementation Roadmap (90 Days)

VAOM implementations follow a five-phase approach designed to deliver measurable results within a single quarter while establishing the controls needed for safe scaling. Phase 1 follows the Delegation Discovery & Design method (Section 5): conduct the Decision Inventory, complete Authority Decomposition for each decision point, select delegation patterns, and validate readiness. Phase 1 deliverables include the completed Decision Inventory, decomposition profiles, pattern assignments, readiness assessment results, and a draft Delegation Authority Matrix.

| Phase | Focus | Key Deliverables | Timeline |
|-------|---|--|-------------|
| 1 | Delegation Discovery & Authority Mapping | Decision Inventory, Authority Decomposition profiles, delegation pattern assignments, Readiness Assessment results, draft Delegation Authority Matrix, stakeholder alignment | Weeks 1-2 |
| 2 | Delegation Design | Delegation Authority Matrix finalized, confidence threshold policy, escalation rules, no-automation zones defined | Weeks 3-4 |
| 3 | Control Alignment | Integration requirements, RBAC mapping, audit log schema, policy library, regulatory control matrix | Weeks 5-7 |
| 4 | Implementation & Wiring | Logging, observability, version control, oversight UIs, human review workflows connected to existing systems | Weeks 8-10 |
| 5 | Pilot & Calibration | Controlled pilot on selected workflow, threshold calibration against real outcomes, authority boundary refinement, audit dry-run | Weeks 11-13 |

Each phase produces evidence artifacts that serve both operational needs and regulatory compliance. The 90-day timeline assumes a single workflow; parallel implementations are possible with additional resources.

Conclusion

AI introduces probabilistic judgment into enterprise workflows. Judgment requires structured authority design. Without it, organizations face a choice between two unsatisfying extremes: blocking AI adoption entirely, or deploying it without adequate controls and hoping the governance gap doesn't surface during an audit or incident.

The Verklöde Agent Operating Model provides a third path: controlled delegation, where automation is bounded by explicit authority, evidenced by immutable audit trails, and accountable to human oversight. VAOM does not require organizations to replace their existing systems, frameworks, or governance structures. It operationalizes them for the age of AI agents.

| *Autonomy, bounded. Decisions, traceable. Accountability, preserved.*

Using VAOM

VAOM is published as an open framework under Creative Commons Attribution 4.0 (CC BY 4.0). Organizations, consultants, and platform teams are free to adopt, adapt, and build on it with attribution to Verkflöde AB.

To apply VAOM to a specific workflow, start with one process. Run the Delegation Discovery & Design method: inventory the decisions, decompose their authority, select patterns, assess readiness. Walk the seven layers. Define what is automated, what is reviewed, what is prohibited. Map the confidence thresholds to your organization's risk appetite. Document the delegation authority for each decision type. The 90-day roadmap in Section 11 provides a starting structure.

The interactive version of this framework, including stakeholder-specific views and the explorable Delegation Discovery method, is integrated into the web version of this document at **verkflode.com/vaom**. The whitepaper PDF is available for download from the same page.

Contributions, feedback, and implementation case studies are welcome at hello@verkflode.com.

Version History

| Version | Date | Changes |
|---------|------|--|
| 1.0 | 2025 | Initial release. Seven-layer architecture, delegation authority matrix, worked example, regulatory alignment, 90-day roadmap. |
| 2.0 | 2026 | Added Delegation Gap scenario, "Why This Is Not Workflow Automation" section, expanded executive summary with scope and architectural neutrality. |
| 2.5 | 2026 | Added "How the Composite Confidence Score Works" (four-dimension scoring with weighted averages and hard floors) and "Managing Foundation Model Drift" (regression testing, threshold recalibration, shadow periods, model registry). |
| 3.0 | 2026 | Added Section 5: Delegation Discovery & Design. Four-stage method (Decision Inventory, Authority Decomposition, Delegation Pattern Selection, Delegation Readiness Assessment). Introduced six named Delegation Patterns and five-dimension Authority Decomposition framework. Updated implementation roadmap and worked example to reference the discovery method. |
| 3.5 | 2026 | Added Section 2: What VAOM Is (and What It Is Not), including credit risk decisioning analogy. Added hidden decisions guidance to Decision Inventory (Section 5.1) with three discovery techniques. Expanded Pattern 6 (Coordinate & Escalate) with multi-agent failure modes: authority conflicts, cascading confidence erosion, and coordination state loss. Expanded Delegation Readiness Condition 5 with guidance on navigating contested ownership. Added Implementation Challenges subsection to composite confidence scoring (Section 8) addressing model calibration, anomaly noise, and rule-model conflicts. Added Calibration Anti-Patterns (five named anti-patterns with symptoms, root causes, and signals). Added three worked examples beyond finance: customer complaint escalation, AML transaction triage, and HR policy violation assessment. |