# Practical Applications of Generative AI in Consulting and Corporate Strategy

## A Governance-First Maturity Ladder for Professional Advisory Work

Alejandro Reynoso

Chief Scientist, DEFI Capital Research
External Lecturer, Judge Business School, University of Cambridge

January 2026

# Disclaimer and Legal Notice

**Educational purpose only.** This book is provided solely for educational and informational purposes. It is not intended to provide, and does not constitute, legal advice, tax advice, accounting advice, investment advice, management consulting advice, or any other form of professional advice. The materials are designed to illustrate governance-first concepts, workflows, and educational examples related to the use of generative artificial intelligence in consulting and corporate strategy. Readers should not rely on this book, or on any example, prompt, workflow, or notebook referenced herein, as a substitute for independent professional judgment or advice tailored to their specific facts and circumstances.

**No professional relationship.** Nothing in this book creates, or should be construed as creating, a client–advisor relationship, consultant–client relationship, fiduciary relationship, or any other professional relationship between the author and the reader. Use of this book does not establish any duty of care, duty of loyalty, or other professional obligation on the part of the author or any affiliated organization.

**No warranties; use at your own risk.** The content is provided "as is," without warranties of any kind, express or implied. The author makes no representations or warranties regarding the accuracy, completeness, reliability, timeliness, or suitability of any information, examples, workflows, prompts, or outputs described in this book. Generative AI systems are probabilistic by nature and may produce incorrect, incomplete, misleading, or fabricated outputs. Any reliance on the material is strictly at the reader's own risk.

**Limitation of liability.** To the maximum extent permitted by applicable law, the author and any affiliated entities shall not be liable for any direct, indirect, incidental, consequential, special, punitive, or exemplary damages arising out of or related to the use of, or inability to use, this book or any associated materials, including but not limited to loss of profits, loss of business, loss of data, reputational harm, regulatory action, or litigation, even if advised of the possibility of such damages.

**Regulatory, legal, and organizational obligations.** Readers remain solely responsible for complying with all applicable laws, regulations, professional standards, contractual obligations, confidentiality requirements, data protection laws, and internal organizational policies. The book does not purport to interpret or replace regulatory guidance, professional standards, or firm-specific governance requirements. Any implementation of ideas discussed herein must be reviewed and approved by appropriate legal, compliance, risk, and management functions.

**Human review required.** All AI-assisted outputs discussed or illustrated in this book are

drafts only and must be reviewed, validated, and approved by qualified human professionals before any reliance, internal decision-making, or external communication. The designation "Not verified" should be assumed unless and until a competent human reviewer has confirmed accuracy and appropriateness.

# Preface

Generative AI is already inside consulting and corporate strategy, whether a firm welcomed it through a formal policy or discovered it living in the quiet spaces between meeting notes, draft decks, internal chats, and the last-minute email to the CFO that must be sent before the taxi reaches the airport. Teams use it to tighten writing, generate first-pass issue trees, produce a list of plausible alternatives, and translate messy workshop transcripts into clean narratives. In-house strategy groups use it to convert a vague executive question into a structured problem statement, to draft an interview guide, or to assemble a decision memo that at least resembles a decision memo. Meanwhile, clients and internal stakeholders increasingly use the same tools to rewrite what you send them, to compress your reasoning into their preferred format, or to produce a counter-memo that looks credible enough to cause trouble. Even when a team tries to avoid AI, it will often encounter AI-generated text indirectly, because it is already embedded in the ecosystem of modern corporate communication.

This book is written for management consultants and in-house strategy professionals who want a clear, non-hyped mental model of what generative AI can do today, what it cannot do, and what must change in professional workflows as capability expands. The goal is not to automate judgment. The goal is disciplined support for structured thinking, scenario analysis, and executive communication. The reader should leave with usable workflows, a clear understanding of limits, and a defensible control framework. The posture is governance-first by design: AI may help you draft and reason, but it may not replace professional responsibility. Every output is treated as a draft, labeled *Verification status: Not verified. Human review required.* unless validated by a human. Facts must be separated from assumptions. Open questions must be made explicit and assigned for verification. If a workflow cannot support reconstructability, it is not ready for real work.

Consulting and corporate strategy have always had a tension between speed and rigor. The modern world rewards speed: boards want answers quickly, CEOs want clarity, and markets punish delay. But rigor is the substrate of trust. When you send a memo to a senior executive, you are not sending prose; you are sending a claim about reality, even if you dress it up as "preliminary." In that context, generative AI is both helpful and hazardous. It is helpful because it can accelerate the production of structure: it can propose a skeleton for a workplan, organize an issue map, suggest a set of plausible drivers, or draft the first version of a storyline. It is hazardous because it can do all of that with a confidence that feels like competence. The central risk is not that AI will write poorly. The central risk is that it will write well enough to pass through your organization as if it

were verified, when it is not.

A helpful way to think about this technology is not as a single tool, but as a ladder of maturity. At Level 1, AI behaves like a chatbot: it is a drafting assistant. You can ask it to rewrite an email, summarize a document, draft an agenda, or produce the first version of a client communication. That is already valuable, and it is already risky, because drafting is where inaccuracies can be smuggled into professional language. At Level 2, AI becomes a reasoner in the practical sense: it can help create issue trees, generate alternatives and tradeoffs, surface assumptions, and design verification plans. At Level 3, AI becomes agentic in workflow form: it can execute a sequence of steps, but only under human checkpoints and explicit stop conditions. At Level 4, the focus shifts from one-off prompts to reusable assets: playbooks, evaluation harnesses, red-team suites, and internal training artifacts that allow safe reuse under time pressure. At Level 5, the organization itself becomes the unit of analysis: governance, supervision, monitoring, approvals, and recordkeeping are embedded into an operating model so that AI support can scale without scaling harm.

The discipline of this book is that every increase in capability introduces new risks, and therefore stronger controls. The capability curve is seductive. It invites you to think, "If the model can draft, it can decide." That is exactly the wrong direction. In professional advisory work, the ability to produce fluent text is not equivalent to the ability to produce defensible conclusions. The only responsible way to use AI at scale is to treat governance as a first-class design constraint, not an afterthought. In a governance-first posture, you do not begin by asking, "What can the model do?" You begin by asking, "What must the firm be able to defend?" Then you work backward: what inputs are permissible, what outputs are allowed, what must be verified, what must be logged, and who owns the final decision.

This is why the book is deliberately didactic. It is written for professionals with minimal AI background, but it does not treat them as naive. Consulting and strategy professionals already understand what disciplined work looks like: assumptions, ranges, scenarios, and the difference between inference and evidence. The problem is not that the profession lacks discipline. The problem is that AI can dissolve discipline by making it easy to skip steps. A well-designed workflow forces you to do the work you would do anyway: clarify the question, identify what must be true, define what you need to know, and create a plan to verify it. AI can help you do those steps faster, but it cannot justify doing them less.

Throughout the chapters, we revisit recurring cases to demonstrate how the maturity ladder changes the workflow. The cases are chosen because they resemble the real work that occupies both consulting teams and in-house strategy groups. They are intentionally familiar: market entry, cost transformation, capital allocation, and operating model redesign. Each case is revisited across levels to show progression. In Level 1, the chatbot helps draft a memo or an email. In Level 2, the reasoner helps generate an issue tree and a verification plan. In Level 3, an agentic workflow executes structured steps with human checkpoints and produces a review packet. In Level 4, the work becomes reusable: templates, tests, and controlled releases ensure that the next team can use the asset without inheriting hidden risks. In Level 5, the organization integrates the workflow into

its operating model so that intake, classification, routing, QA, approvals, and recordkeeping are standardized and auditable.

The recurring cases are also a pedagogical choice. Many AI materials feel like a collection of tricks. They show a prompt that produces an impressive output, then move on. That style is not useful in professional practice because it fails to teach what matters: process. Consulting is not a single deliverable; it is a sequence of decisions under uncertainty, with accountability. The purpose of this book is therefore not to provide "magic prompts," but to provide workflows that survive stress: incomplete information, conflicting stakeholder demands, time pressure, and the temptation to treat fluent output as truth.

The most important concept you will see repeated is the separation of facts and assumptions. Consulting teams are trained to use assumptions. In fact, assumptions are often the only way to make progress when the data is incomplete. The problem arises when assumptions are implicit, unowned, and untested. AI makes that problem worse because it will happily fill gaps with plausible content unless you constrain it. The book therefore insists on a simple standard: facts must be provided or verified; assumptions must be stated and owned; open questions must be listed and prioritized; and every output must include a verification plan. If you adopt only one practice from this book, adopt that one. It will prevent most of the failure modes that make AI dangerous in professional work.

Another core concept is decision laundering: the practice of presenting an AI-generated conclusion as if it were a neutral authority, thereby reducing perceived human responsibility. This can happen subtly. A team may produce a recommendation and then say, "The model suggests..." as if that phrase reduces the burden of justification. It does not. AI does not have standing in governance. It is not accountable. It cannot be cross-examined. It does not bear reputational or legal consequences. The only accountable actor is the human professional and the organization. This book therefore prohibits "autonomous strategy recommendations." The model may draft, structure, enumerate, and propose alternatives, but the judgment must remain human, and the rationale must remain documented.

That constraint also shapes how we treat numbers, benchmarks, and citations. This book does not allow fabricated data. It does not allow invented benchmarks. It does not allow invented citations. If a statistic is needed, it must be supplied as an input or explicitly marked as an assumption to verify. If a reference is claimed, it must be verified by a human before it is treated as real. This is not merely a preference; it is a governance requirement. In consulting, numbers and citations are power. They anchor narratives. If you allow a model to invent them, you create a system that can generate false confidence at scale. Governance-first means you would rather be slower than wrong in ways that are undetectable.

The maturity ladder also forces a distinction between individual use and institutional use. A single person can often use AI safely through personal discipline: careful prompting, cautious tone, and manual verification. The moment you scale to a team, discipline becomes inconsistent. People copy and paste outputs. They reuse drafts without knowing their provenance. They inherit

assumptions without seeing them. At scale, safety is not a matter of individual virtue; it is a matter of system design. That is why Levels 4 and 5 are about assets and organizations, not prompts. A reusable asset must embed scope boundaries, structured outputs, evaluation tests, and versioning. An organization must embed approvals, recordkeeping, monitoring, and retirement mechanisms. In this framing, governance is what allows reuse to scale without scaling harm.

The companion Google Colab notebooks are a key part of that system. Each chapter has a notebook that implements the chapter's maturity level. These notebooks are educational, but they are designed to produce the kinds of artifacts that real teams need: run manifests, prompt logs, risk logs, and deliverables bundles. Every run produces a saved trail: what inputs were provided, what prompts were used, what outputs were generated, what risks were flagged, and what files were produced. The notebook outputs are not "the answer." They are drafts and evidence. They are designed to be reviewed, critiqued, and improved. The discipline is that every output is labeled *Verification status: Not verified. Human review required.* by default until a human validates it.

This audit artifact posture is not cosmetic. It is the foundation of defensibility. In real professional practice, you should assume that any important output may need to be reconstructed later: for internal review, for client challenge, for regulatory inquiry, for litigation, or simply for learning. If you cannot reconstruct the work, you cannot defend it. AI-generated content increases the need for reconstructability because it introduces a probabilistic component to drafting: the same prompt may not produce the same output, and the model may behave differently over time. The only way to manage that is to log what happened, version what you used, and maintain a record of human review and sign-off.

Each chapter also contains explicit "CAN" and "CAN'T" statements. Those are not motivational slogans; they are control boundaries. For example, a Level 1 chatbot can draft a memo, but it cannot validate facts. A Level 2 reasoner can propose an issue tree, but it cannot decide which branch is true. A Level 3 workflow can execute steps, but it must stop if hinge facts are missing. A Level 4 asset can be reused, but only if it has evaluations and controlled release. A Level 5 organization can scale, but only if it can monitor, audit, and retire assets. These boundaries are the practical difference between responsible use and performative use.

The book's tone is intentionally conservative in one specific way: it never treats AI output as inherently reliable. That is not pessimism. It is a governance stance. In professional work, the cost of persuasive error is high, and the cost of undetected error is higher. AI is uniquely capable of producing undetected error because it can produce language that looks correct even when it is not. Your job is to design workflows that surface uncertainty rather than hide it. This is why we repeatedly require the model to list assumptions, open questions, and verification steps. The model is not being asked to be omniscient. It is being asked to be transparent.

If you are a consultant, you will recognize the meta-problem: clients often want crisp answers, while the truth is conditional. AI can make it tempting to give them what they want: crispness. A governance-first posture insists that crispness is earned, not generated. It comes from verification, not from phrasing. If you are an in-house strategy leader, you will recognize a similar dynamic:

executives want decisions, while the organization has incomplete data. AI can accelerate the organization's ability to frame the decision, but it cannot fix the underlying incompleteness. The professional discipline is to use AI to make the incompleteness visible and manageable, rather than to use AI to obscure it.

A brief note on what this book is not. It is not a catalog of vendor features. It is not a benchmark report. It is not an argument that AI will replace consultants. It is not an attempt to turn consulting into automation. It is an operating manual for a specific thesis: AI can support structured thinking and drafting at scale only if governance is designed first. If you are looking for a hype-driven narrative about disruption, you will find this book boring. That is by design. Real professional practice is often boring when done correctly. It is repetitive, disciplined, and documented. The book's claim is that AI can make that boring discipline faster and more consistent, but only if you refuse to let fluent output substitute for evidence.

Finally, a word about adoption. Most organizations do not fail because they chose the wrong model. They fail because they adopted AI without changing their process. They treat AI like a faster intern, then act surprised when errors appear. They deploy prompts without controls, then act surprised when outputs drift. They encourage experimentation without logging, then act surprised when they cannot reconstruct what happened. If you want to avoid that path, adopt the ladder deliberately. Start with Level 1 and build habits. Introduce Level 2 structures so that reasoning becomes inspectable. Add Level 3 workflows only when checkpoints and logs are real. Build Level 4 assets only when you can version and test them. Move to Level 5 only when the organization is ready to own the operating model.

If you read this preface as a warning label, you will miss its optimism. The optimism is that consulting and corporate strategy already have the intellectual tools needed to use AI responsibly: structured thinking, explicit assumptions, scenario planning, and professional accountability. AI does not need to invent a new discipline. It needs to be placed inside the existing discipline, with controls that acknowledge its failure modes. If you do that, you can get the real benefits: faster drafting, better structured reasoning, improved consistency, and more defensible executive communication. You will also get something rarer: a practice where the use of AI strengthens trust rather than weakening it.

The chapters that follow are therefore written as a progression of capability and control. Each chapter provides practical workflows and copy/paste prompt patterns, but also insists on governance artifacts, stage boundaries, and review discipline. The companion notebooks make that discipline executable. The destination is not autonomy. The destination is a professional operating system where AI accelerates work without eroding accountability. That is the maturity ladder. That is the governance-first promise. And that is what this book aims to deliver.

# How to Use This Book

This book is designed to be used the way consulting work is actually produced: iteratively, under time pressure, with incomplete information, and with real consequences for clarity, accuracy, and trust. It is not a showcase of clever prompts. It is a governance-first operating manual for using generative AI without surrendering professional responsibility. The aim is to help you become faster at the parts of the job that should be fast (drafting, structuring, summarizing, producing first-pass alternatives), while becoming stricter about the parts of the job that must be strict (verification, assumptions, decision rationale, confidentiality, and sign-off). If you approach it that way, each chapter becomes a reusable workflow you can apply tomorrow, with guardrails that scale as capability increases.

Read the book in order. The five chapters are a maturity ladder, not a menu. Level 1 shows how to use AI as a drafting assistant without importing unverified claims into client deliverables. Level 2 introduces explicit reasoning structures: issue trees, hypotheses, alternatives and tradeoffs, assumption registers, and verification planning. Level 3 adds multi-step workflows with human checkpoints, separation of duties, and audit-ready logs. Level 4 turns workflows into reusable internal assets that can be versioned, tested, and safely distributed across teams. Level 5 lifts the entire system to the organizational level, where governance is not a guideline but an operating model: intake controls, classification, routing, QA, approvals, and recordkeeping that allow reconstruction after the fact. Skipping ahead is tempting, but it is usually how teams end up with impressive demos and fragile practice.

Use a simple rhythm: one chapter per week (or per two weeks), paired with its companion Colab notebook. Read the chapter first, then run the notebook. Do not treat notebook outputs as deliverables; treat them as drafts and evidence. The notebooks are intentionally designed to produce artifacts that look like consulting work product, but their real purpose is to train disciplined process. In every run you will see the same governance posture repeated: outputs are labeled *Verification status: Not verified. Human review required.* by default; facts and assumptions are separated; open questions are explicit; and the run produces a saved bundle of logs and deliverables. This repetition is not redundant. It is the point. In real client work, the team that wins is rarely the one with the most creative prompt; it is the one with the most reliable workflow under pressure.

A useful mental model is to treat AI as a junior analyst who writes quickly, speaks confidently, and does not know what it does not know. That analyst can be extremely helpful if you give precise instructions, constrain scope, and require evidence. That analyst can be dangerous if you allow it

to improvise sources, "fill in" missing facts, or translate uncertainty into fluent certainty. Your job is to manage the interface so that speed does not become a substitute for truth. Throughout the book, whenever you see a rule that feels conservative, assume it was written in response to a failure mode that looks professional in slide form and catastrophic in governance form.

The first habit to adopt is separation: facts, assumptions, analysis, and verification. If the input does not contain a fact, the output must not pretend it does. If the team needs a benchmark, a statistic, or a quote, the workflow must either (a) require the user to provide it, or (b) label it as an assumption to be verified and assign an owner to verify it. The book treats "facts" as items that are either provided directly by the user (for example, numbers from the client's ERP extract, or statements from an internal memo) or verified by the team against a trusted source. Everything else is an assumption. This is not pedantry; it is the only way to keep AI-assisted outputs from quietly drifting into decision-grade language without decision-grade evidence.

The second habit is scoping: define what the model is allowed to do, and what it is not. In consulting, ambiguity expands by default. AI amplifies that. Each chapter therefore emphasizes scope boundaries: what inputs are in-bounds; what is out-of-bounds; what the output is allowed to claim; and what the output must explicitly refuse to do. When you use the prompts and workflows in this book, do not weaken those boundaries to "make it work." Instead, treat boundary failures as signals that you need better inputs, clearer objectives, or a different workflow stage. The fastest way to degrade governance is to teach the system that it should always produce an answer, even when the right answer is "we cannot conclude this from the information provided."

The third habit is verification planning. AI is often useful not because it knows, but because it can help you identify what you need to know. At Levels 1 and 2, verification planning shows up as checklists and "questions to verify." At Level 3 and above, it becomes stage gates: the workflow is allowed to draft, but it is not allowed to finalize until specific hinge facts are confirmed. Verification planning should be operational, not rhetorical. That means owners, timestamps, and acceptance criteria. It also means making "stop conditions" normal: if the client has not provided revenue by segment, the market sizing portion of the storyline is blocked; if legal has not approved a claim, it cannot appear in an executive memo; if finance cannot validate a number, it cannot enter a model.

The fourth habit is recordkeeping. Consulting work is increasingly judged by its reconstructability: can you explain what you did, what you relied on, and why the conclusion follows from the evidence? The notebooks reinforce this by generating run artifacts each time: a run manifest (what ran, when, and under what configuration), a prompt log (what was asked and what was produced), a risk log (what could go wrong and what controls were triggered), and a deliverables bundle (what the workflow produced). In practice, you should treat these artifacts as the skeleton of a defensible workfile. They do not replace project documentation; they raise the baseline. They also make supervision real: a reviewer can examine what happened, rather than reviewing only polished prose.

The fifth habit is role clarity. Consulting is a team sport with implied separation of duties: an analyst drafts, a manager challenges, a partner owns the client narrative. AI can blur those boundaries if it becomes a single "black box" that produces everything. The maturity ladder reverses

that temptation. At Level 1, the human is the author and the model is a drafting assistant. At Level 2, the model helps structure reasoning, but it cannot claim conclusions. At Level 3, workflows can be multi-step, but the human checkpoints are explicit and mandatory. At Level 4, assets are reusable only if they embed controls and evaluation. At Level 5, the organization owns the operating model, including approvals, monitoring, and retirement of assets. Use this progression as a management tool: do not let AI become a substitute for managerial challenge, and do not let "the model said" become a rationale.

To make the book practical, each chapter includes copy/paste prompt patterns and exercises. Use them as templates, not as scripts. Replace placeholders with your actual inputs. Keep the constraints that force discipline: the request for structured outputs, the insistence on "Not verified," and the requirement to list assumptions and open questions. If you are running these prompts in a live engagement, add one more constraint: confidentiality. Minimize what you paste. Redact client identifiers. Avoid sensitive details unless you are in an approved environment. In a governance-first posture, it is always acceptable to provide less and ask the model for structure rather than specifics. A well-designed workflow can produce useful scaffolding with minimal data.

Finally, understand what this book is not trying to do. It is not trying to convince you that AI will replace consultants or strategy teams. It is not trying to provide a universal library of benchmarks or a catalog of "best practices." It is not trying to automate judgment or outsource responsibility. It is trying to help you design a disciplined interface between human professionals and probabilistic text generation, so that speed serves rigor rather than undermining it. If you treat the maturity ladder as a governance ladder, you will get the intended benefit: faster drafting, clearer reasoning structures, better verification habits, and more defensible executive communication.

If you want a simple scorecard for whether you are using the book correctly, use this: after running a workflow, can you point to (1) the facts you relied on, (2) the assumptions you made, (3) the open questions you still need answered, (4) the tests you plan to run, and (5) the human owner who will sign off? If the answer is yes, you are building professional-grade practice. If the answer is no, the output may look impressive, but it is not governed. The ladder is not about capability. It is about defensibility. Use it accordingly.

# Chapter 1

# Chatbots

**Abstract.** This chapter introduces Level 1 of the governance-first maturity ladder for AI consulting and corporate strategy: the disciplined use of chatbots for drafting professional artifacts. At this level, artificial intelligence is treated neither as an analyst nor as a source of truth, but as a constrained drafting engine that accelerates formatting, language, and structural clarity while leaving judgment, verification, and accountability entirely with the human professional.

The chapter establishes a clear mental model for safe Level 1 usage. Chatbots can convert rough notes into coherent emails, memos, decision records, and outlines, but they do not "know" facts, validate assumptions, or check reality. Because their outputs are often fluent and persuasive, the primary risk at this level is not technical error but managerial error: unverified drafts being mistaken for reliable analysis. The governance posture therefore emphasizes explicit separation of facts provided, assumptions made, and open questions requiring verification.

Readers are given a minimum control framework suitable for consulting teams and corporate strategy functions. This includes a standard drafting wrapper, redaction and minimum-necessary input rules, versioning discipline, and a mandatory "Not verified" status for all AI-generated drafts. The concept of a saved "Draft Bundle" is introduced to ensure traceability, reviewability, and audit readiness from the first use of AI.

By the end of the chapter, the reader will be able to use chatbots productively without over-delegating responsibility, produce cleaner and more transparent drafts, and establish habits that scale safely into higher levels of AI capability.

---

**Artifact (Save This)**

**Scope disclaimer** This chapter is an educational governance-first guide for management consulting and corporate strategy work. Outputs produced with AI are drafts only and must be reviewed by a qualified professional. The model may be persuasive and wrong. Do not treat outputs as factual, complete, or client-ready without verification and sign-off.

**Level alignment.** Level 1 uses a chatbot to accelerate drafting and formatting. No autonomous agents, no tool use, no independent research is assumed. Human judgment remains fully responsible.

## 1.1 Chapter overview: Level 1 in the maturity ladder

Level 1 is where most organizations will begin, and where most organizations will remain for longer than they expect. That is not a failure of ambition. It is an accurate reflection of the fact that the first meaningful value of generative AI in consulting and corporate strategy is not "strategy on autopilot," but disciplined drafting: turning partial notes into coherent artifacts, imposing structure where humans are rushed, and making communications clearer without claiming that anything has been verified. Level 1 is therefore the entry point of a governance-first ladder: the moment when an organization chooses to treat AI as a productivity tool *and* as a new source of operational risk. This chapter frames Level 1 as a managerial decision, not a technical novelty. The question is not whether the model can write. The question is whether the team can use it without confusing fluency for truth, and without weakening the professional discipline that makes consulting and strategy credible.

The maturity ladder matters because it provides a shared language for capability, risk, and control. Without a ladder, AI adoption becomes a chaotic mix of individual experimentation, inconsistent practices, and untraceable outcomes. One team uses the model for harmless rephrasing; another uses it to generate "market benchmarks" that were never provided; a third pastes confidential client data into an unapproved interface because it "worked last time." A maturity ladder is a governance instrument: it allows leaders to define what is permitted at each stage, what must be logged, what must be reviewed, and what is explicitly out of scope. It also creates a practical training path. Consultants and strategy professionals do not need to become machine learning engineers to use AI safely, but they do need a disciplined mental model: what the tool is, what it is not, and what control measures are proportionate to the capability being used.

Level 1 should be understood as the foundation. It is the smallest set of practices that allows a team to benefit from chatbots without introducing unacceptable risks. It does not attempt to solve every problem. It does not promise "insights." It does not replace analysis. Instead, it establishes habits that are easy to follow, easy to audit, and easy to scale. If those habits are not present at Level 1, every higher level becomes dangerous, because higher levels amplify both the upside and the failure modes. If a team cannot reliably separate facts from assumptions in a one-page email draft, it is not ready to let AI generate multi-step deliverables, orchestrate workflows, or ship reusable assets.

At the same time, Level 1 is not trivial. Drafting is where most professional work begins and where many failures originate. Strategy failures often start as narrative failures: executives being briefed with false certainty, stakeholders being misaligned because an email concealed uncertainty, or teams committing to plans whose assumptions were never written down. In consulting, the earliest drafts of materials often shape the entire engagement: how the problem is framed, what options are deemed plausible, what risks are treated as "known," and what evidence is considered sufficient. A chatbot, used carelessly, can accelerate the wrong direction. Used with discipline, it can strengthen the communication layer of professional judgment: clearer drafts, cleaner structures, better questions, and fewer hidden assumptions.

This chapter therefore begins by positioning Level 1 in the ladder, and by clarifying a central principle that will repeat throughout this book: every increase in AI capability introduces new risks and therefore requires stronger controls. The ladder is not a marketing slogan. It is a governance framework. It is not "Level 5 is better"; it is "Level 5 is more powerful and therefore more dangerous." The responsible posture is not to race upward, but to build competence and controls at each step, with explicit scope boundaries and auditable practices.

### 1.1.1 The five levels (why this ladder exists)

The five levels are a pragmatic abstraction of how AI capability is typically applied in professional contexts. They are not a technical taxonomy of model architectures. They are a managerial map: what the tool is being used to do, how much autonomy it is given, and what governance must exist to make that usage defensible. Each level can be adopted in a conservative or aggressive way, but the defining difference is the *degree of delegation* and the *degree of workflow complexity.* In consulting and corporate strategy, the key question is always: who owns the judgment, and how is that ownership evidenced?

**Level 1: Chatbots (drafting and formatting).** At Level 1, the model is used as a drafting assistant. Inputs are provided by the human. Outputs are drafts to be reviewed by the human. The model is not asked to discover facts, to "research," or to validate truth. It can reorganize, rewrite, summarize, and structure. It can propose templates and agendas. It can generate alternative phrasings and tones. But it must not be treated as an authoritative source. Governance at this level is largely about preventing persuasive error and accidental leakage: redaction discipline, schema discipline (facts vs assumptions vs open questions), and a review gate before anything is sent externally. The Level 1 posture is straightforward: the model accelerates writing, while humans own reality.

**Level 2: Reasoners (structured thinking and explicit assumptions).** At Level 2, the model is used to support structured reasoning: issue trees, option lists, tradeoffs, scenario narratives, and explicit assumption registers. The model still does not "know" facts, but it can help structure the problem and expose gaps. The central value is not writing quality, but thinking quality: better decomposition, clearer alternatives, more explicit uncertainty. Governance must strengthen because the model is now shaping analytic framing, not merely phrasing. Controls include assumption registers that block downstream steps when unverified, reasoning templates that force transparency, and more rigorous reviewer expectations. Level 2 is where teams often feel a temptation to treat the model as an "analyst." The governance-first framing insists otherwise: it is a reasoning scaffold, not a truth engine.

**Level 3: Agents (multi-step workflows with human checkpoints).** At Level 3, AI is embedded in a multi-step workflow. Instead of one prompt and one output, there is an orchestrated sequence: intake, classification, drafting, quality checks, red-team pass, escalation of uncertainty, and packaging into deliverables. Some steps may be automated, but the workflow contains explicit human checkpoints. Governance becomes process governance: separation of duties, immutable logs,

defined owners, and error handling. The risk is no longer just hallucination; it is workflow failure, silent drift, and unobserved compounding errors across steps. At Level 3, defensibility requires audit artifacts: run manifests, prompt logs, risk logs, and reviewer sign-offs.

**Level 4: Innovators (reusable assets and controlled releases).** At Level 4, the organization begins to package AI workflows into reusable assets: standardized prompt libraries, internal playbooks, tested templates, and supervision frameworks that can be applied across teams. The key shift is from "use" to "productization." Governance becomes akin to internal product governance: version control, testing, release notes, scope constraints, and a controlled change process. The risks include brittle reuse, hidden assumptions baked into templates, and uncontrolled proliferation. Controls include evaluation harnesses, regression tests, and documented intended use. Level 4 is where firms begin to create genuine IP, but also where they can inadvertently create scalable failure.

**Level 5: Organizations (firm-level governance and auditability).** At Level 5, AI is treated as a governed organizational capability. There is a system of record. There are policy-defined intake processes, data controls, role-based permissions, training requirements, QA standards, and auditability by design. Model usage is not merely "allowed"; it is managed, monitored, and continuously improved. The organization can demonstrate how AI outputs are produced, reviewed, and approved, and can defend those processes to internal risk teams, clients, regulators, and auditors. Level 5 is not about maximum autonomy. It is about maximum accountability. It may enable advanced use cases, but it demands disciplined controls that match enterprise risk expectations.

A ladder like this exists because organizations routinely confuse *what a tool can generate* with *what a tool can be trusted to decide.* In consulting and strategy, the work product is not just text. It is the chain of reasoning, the evidence base, the accountability for assumptions, and the professional responsibility for recommendations. AI systems can produce plausible language quickly. That plausibility creates a new failure mode: the output looks credible enough to be believed. A maturity ladder is a counterweight. It forces teams to connect capability to control: if you want more autonomy, you must accept more governance.

The ladder also exists because it allows staged adoption. Many organizations do not need Level 5 to capture most of the value of AI. In fact, many organizations should deliberately remain at Level 1 or Level 2 for long periods, because their risk appetite, confidentiality constraints, and operating model maturity do not support agentic workflows or reusable assets. The ladder gives permission to be conservative. It also makes conservative adoption legitimate: "We are at Level 1 by design, because this is what we can govern today."

Another reason the ladder exists is training. Consulting firms and strategy functions are often comprised of high-performing generalists. They adopt tools quickly, but they do not always adopt controls. A ladder becomes a curriculum: each level introduces new skills and new obligations. The ladder teaches that governance is not a compliance afterthought. It is the core professional discipline that allows AI to be used without undermining trust.

Finally, the ladder exists because it encourages explicit boundaries. A common organizational failure is to adopt AI in a vague way: "Use it when helpful." That instruction is an invitation to

inconsistency. One person uses it to rewrite grammar; another uses it to generate competitive intelligence. A ladder replaces vagueness with defined scope: "At Level 1, you may do X and you may not do Y." This is essential in consulting and strategy because the cost of a mistake is often reputational and irreversible. A single unverified "industry benchmark" slipping into a client deck can damage credibility. A single confidential detail pasted into an unapproved interface can violate policy. In these contexts, adoption without boundaries is adoption without governance.

### 1.1.2   What changes at Level 1 (and what does not)

Level 1 changes one thing: the speed and ease with which professionals can produce readable drafts. It does not change the underlying responsibility structure of consulting and corporate strategy work. The human remains the accountable author. The team remains responsible for the truth of what is said. The organization remains responsible for confidentiality and data governance. The model does not become a researcher, a fact-checker, or an expert witness. Level 1 is a drafting accelerator. It is powerful precisely because drafting is pervasive. It is risky precisely because drafting can masquerade as analysis.

To adopt Level 1 responsibly, it is essential to specify what changes and what does not in operational terms.

**What changes: the drafting workflow becomes faster and more structured.** Most consulting and strategy artifacts begin as messy inputs: bullet notes from a meeting, partial thoughts from a partner, fragments of analysis from multiple workstreams, and a timeline that does not permit careful writing. Level 1 enables teams to transform those messy inputs into coherent drafts with consistent structure. This includes:

- **Rapid translation from notes to narrative.** A set of disorganized bullets can be reshaped into a memo that has a clear context, objective, key messages, and next steps.
- **Standardization of format.** Teams can enforce consistent memo structures, email templates, decision records, and meeting summaries, reducing variation and omissions.
- **Tone control.** A draft can be rendered in multiple professional voices (executive concise, diplomatic, assertive) without changing content.
- **Compression and expansion.** A long draft can be compressed into an executive summary, or a short draft can be expanded into a fuller narrative, as long as the facts are controlled.
- **Checklist generation.** From a described situation, the model can produce a checklist of questions to ask, risks to consider, or stakeholders to engage, serving as a disciplined prompt for the team's own thinking.

These changes are valuable because they reduce time spent on mechanical writing and increase time available for judgment and verification. They can also reduce the risk of accidental omissions. Many professional failures occur because something obvious was not stated: an assumption was not written down, a risk was not surfaced, a dependency was not listed, or an open question was not escalated. A Level 1 drafting wrapper can force these elements into the output format.

**What does not change: the model does not add verified knowledge.** The most important non-change is epistemic. Level 1 does not create truth. A chatbot may generate statements that sound like facts, but unless those facts were provided by the user and verified independently, they are not facts. The model does not have a duty of care. It does not bear reputational risk. It does not suffer consequences for being wrong. It is therefore a mistake to treat it as an analyst who "did the research." In governance-first terms, the model has no standing to introduce new factual claims into a deliverable.

This non-change has a practical implication: at Level 1, you must adopt a "no new facts" rule. The model may reorganize facts you provide, but it must not invent or import facts. If information is missing, it must be recorded as an open question. If a benchmark is needed, it must be labeled as required input rather than guessed. This is not merely an ethical preference; it is a control that prevents persuasive error.

**What does not change: verification remains a human process.** Consulting and corporate strategy are credibility businesses. Verification is not optional. It may be lightweight for some artifacts (a follow-up email), and rigorous for others (a board memo), but it is always necessary. Level 1 outputs should therefore carry an explicit status: `"Not verified"`. This is not pessimism; it is accurate labeling. It prevents internal recipients from treating drafts as final. It also creates discipline. If the output is labeled not verified, the next step is naturally: who will verify, what must be checked, and by when?

**What does not change: accountability and authorship remain human.** A common failure in early AI adoption is "authorship diffusion." The draft came from the model, so no one feels fully responsible for it. In professional work, that is unacceptable. The governance-first posture requires named ownership: who requested the draft, who reviewed it, who approved it, and what evidence exists of that review. Even if the artifact is informal, such as an email, the discipline should exist: a reviewer note, a version history, or at minimum an explicit acknowledgment that the sender has validated the content.

**What changes: teams can create internal drafting IP earlier than expected.** Another subtle change at Level 1 is that teams can begin to standardize their communication patterns. Many consulting teams already have unofficial templates. Level 1 allows those templates to be formalized: a standard memo wrapper, a standard meeting summary format, a standard risk section, a standard "ask" section. These patterns can become internal IP, but at Level 1 they must remain carefully scoped: templates for structure, not templates that smuggle in conclusions. The value is consistency. The risk is hidden assumptions being repeated at scale.

**What does not change: confidentiality requirements become stricter, not looser.** A paradox of Level 1 is that because it feels harmless, teams may be tempted to paste more sensitive data into the tool. Governance-first adoption insists on the opposite: minimum-necessary input. Level 1 is about drafting, so it rarely requires raw confidential material. Names can be replaced with placeholders. Deal terms can be generalized. Identifiers can be removed. The guiding principle is simple: if the model does not need it to produce the draft structure, do not provide it. Redaction

is not an afterthought. It is a prerequisite.

**What changes: managers must now manage AI usage as process, not as individual preference.** At Level 1, leadership's job is not to teach "prompt tricks." It is to define safe usage patterns and enforce them. That means specifying:

- Which deliverable types are permitted (emails, memos, agendas, storylines).
- Which deliverable types are prohibited (final recommendations, benchmark-heavy decks).
- The required output schema (facts, assumptions, open questions, draft).
- The requirement to label outputs as not verified.
- The requirement to save a minimal draft bundle for important deliverables.

This is an operating model decision. It is also a cultural decision. If leaders treat AI as a private productivity hack, they will get inconsistent, untraceable usage. If leaders treat AI as a governed drafting method, they will get repeatable value and defensible outcomes.

**What does not change: professional judgment remains the core of consulting and strategy.** It is worth stating plainly: the purpose of consulting is not to produce text; it is to deliver judgment, supported by evidence, aligned with client context, and communicated persuasively but responsibly. Level 1 can improve the communication layer, but it cannot supply the judgment. In fact, if the team is not careful, Level 1 can degrade judgment by making weak thinking look strong. The governance wrapper therefore exists to preserve judgment: by forcing explicit assumptions, by surfacing open questions, and by preventing invented facts.

To make this concrete, consider a common situation. A partner asks for a one-page memo to send to a CFO about a proposed cost transformation approach. The analyst has partial notes and a tight deadline. A Level 1 chatbot can produce a polished memo quickly. But if the memo includes invented cost savings percentages or implied timelines that were never agreed, it can mislead the CFO and create commitments the team cannot support. The correct Level 1 behavior is different: the memo should clearly separate the facts provided (scope, timeline constraints, known cost categories) from assumptions (estimated savings ranges, required approvals) and open questions (data availability, labor constraints, implementation risks). The draft can still be persuasive. But it must be honest about what is known and unknown.

This is why the ladder begins at Level 1 with discipline. The first lesson of AI in consulting is not "how to prompt." It is "how to maintain professional integrity when language becomes cheap." In a world where a model can produce a plausible paragraph in seconds, the scarce resource is no longer writing time. The scarce resource is verification, judgment, and accountability. Level 1 is the moment when teams must choose to protect that scarce resource rather than accidentally dilute it.

In the chapters that follow, the ladder will add additional capabilities: structured reasoning, multi-step workflows, reusable assets, and organizational systems. Each capability will come with new risks and new controls. But the foundation remains the same: facts are not assumptions; outputs are drafts until verified; and accountability cannot be delegated to a model. If those principles hold at Level 1, the organization can safely capture early value and build toward more advanced maturity. If those principles do not hold at Level 1, every higher level becomes a scaling

mechanism for error.

This chapter overview therefore sets the tone for the rest of the book: governance first, capability second. The objective is not to make AI impressive. The objective is to make AI usable in professional practice without eroding trust. Level 1 is where that trust is either strengthened through discipline, or weakened through shortcuts. The ladder exists to ensure the former, and to prevent the latter.

## 1.2 Mental model: what a chatbot is in consulting work

If Level 1 is the entry point of the maturity ladder, the mental model is the guardrail that keeps the entry point from becoming an accident. Most AI failures in consulting and corporate strategy do not begin with malicious intent or obvious negligence. They begin with a subtle mismatch between what a chatbot *appears* to be doing and what it is *actually* doing. The model produces language with confidence, coherence, and professional cadence. Humans are trained to associate those qualities with competence. In a consulting setting, where polished writing is often the wrapper around real analysis, it is dangerously easy to mistake the wrapper for the substance.

A governance-first mental model solves this by forcing two disciplines at once. First, it defines the chatbot in operational terms: what kind of machine it is, what kind of outputs it produces, and what it cannot do. Second, it defines the human obligations that must remain intact regardless of how fluent the output looks. This chapter does not ask consultants to become AI specialists. It asks them to become epistemically strict: to treat the model as a drafting device that mirrors patterns, not a colleague that verifies truth.

A useful way to think about Level 1 is that it changes the cost structure of communication. Historically, writing a clean memo or email took time. That time constraint functioned as an accidental governance mechanism: when writing is expensive, people pause, reconsider, and refine. When writing becomes cheap, the pause disappears. The risk is not that the model writes poorly; the risk is that the organization ships drafts that feel finished. The mental model must therefore replace the old friction with explicit controls: structure, labeling, and review.

### 1.2.1 Useful abstraction

The safest abstraction for a Level 1 chatbot in consulting work is this: *a drafting engine with pattern recall that functions as a structured mirror of your inputs.* The two parts matter.

**Drafting engine.** The model is best treated as a tool that transforms inputs into a formatted output. It is not a knowledge authority in your context. It is not a fact-checker. It is not a due diligence team. It is a system that can take raw notes and produce a coherent first draft. In consulting, this is already valuable because many deliverables are constrained by time and by the need for consistent structure: meeting summaries, stakeholder updates, decision memos, workshop agendas, interview guides, and slide storylines. If the model can do the first pass of organization and phrasing, the team can spend its scarce time on the higher-value work: verifying facts, checking logic, aligning stakeholders, and refining recommendations.

Calling it a drafting engine is also a governance choice. It sets an explicit ceiling: the model may accelerate expression, but it may not become the source of substantive claims. This ceiling prevents the most common failure mode: invented facts slipping into a deliverable because they were expressed fluently.

**Pattern recall.** The model's strength is its ability to produce plausible patterns of language and structure based on what it has learned. In practice, this means it is good at generating a "consulting-shaped" memo, a "CFO-shaped" email, or a "board-shaped" outline. It knows how those documents tend to look. It can imitate the conventions: headings, bullets, executive summaries, risk sections, next steps. That is exactly why it is helpful in Level 1. It imposes familiar form quickly.

But pattern recall is also why it is dangerous. A pattern engine can produce a convincingly complete artifact even when critical information is missing. In a human workflow, missing information often reveals itself as awkward gaps: a paragraph you cannot write because you do not know the answer, a section you cannot complete without data, a sentence that forces you to pause and ask someone a question. The model does not experience that pause. It fills the gap with plausible completion. The mental model must therefore train you to expect gap-filling and to preempt it by forcing the model into a constrained output schema that makes gaps explicit.

**Structured mirror.** The phrase "structured mirror" captures the most valuable use of Level 1 beyond mere drafting. When you provide the model with your notes and ask it to produce a structured output, you are effectively asking it to reflect your thinking back to you in a cleaner form. This reflection can reveal problems that were hidden in the original mess: contradictions, missing decisions, unclear ownership, ambiguous asks, and unstated assumptions. The model is not discovering truth; it is reorganizing your input in a way that makes truth claims and uncertainty more visible.

In consulting work, this mirror function is practical. After a client call, your notes may contain fragments: "CFO worried about timeline," "Ops says data not ready," "CEO wants quick win," "Need baseline," "Procurement." A structured mirror can transform these into a meeting summary with explicit decisions, action items, owners, due dates, and open questions. When the open questions are written down, you can route them to the right people. When assumptions are forced into a list, you can challenge them early. When "facts provided" are listed, you can see whether you have enough to send anything.

This is the key: the model's output is only as trustworthy as the inputs and the constraints you impose. The structured mirror becomes useful when you use it to make your own work more transparent to yourself and to your team. It becomes dangerous when you treat it as independent analysis.

To operationalize the useful abstraction, Level 1 work should be framed with a handful of consistent instructions:

- **You are drafting only.** The model is not allowed to invent facts or imply verification.
- **Use only the facts I provide.** If data is missing, list it under open questions.
- **Separate facts, assumptions, and open questions.** Do not blend them into prose.

- **Label the output as not verified.** The output is a draft for review.

These instructions are not "prompt hacks." They are governance controls expressed in language the tool can follow.

### 1.2.2 Dangerous misconception

The most dangerous misconception in consulting work is not that the model sometimes makes mistakes. Everyone can accept that. The dangerous misconception is epistemic: that the model *knows* something about the situation that you did not tell it, or that it *checked* something you did not check, or that it is *authoritative* because it sounds authoritative.

In consulting and corporate strategy, authority is a social construct built on credibility: evidence, track record, professional accountability, and explicit ownership of claims. A chatbot has none of these. It can produce language that mimics authority, but it cannot bear responsibility. The governance-first approach therefore treats "model authority" as a category error.

There are several variants of this misconception, and each has a characteristic failure mode.

**Misconception 1: "The model knows."** This appears when teams assume that the model can supply industry facts, benchmarks, competitor moves, or regulatory constraints without being given those facts. In reality, any such output is at best generic and at worst fabricated. The failure mode is "false precision": a cost reduction percentage, a market growth rate, or a timeline that looks realistic but is unsupported. In consulting, false precision is poisonous. It travels quickly through decks and memos, and by the time it is challenged, it has already shaped decisions.

**Misconception 2: "The model checked."** This appears when teams treat the model as if it validated logic, cross-checked numbers, or ensured internal consistency. A model can help you *notice* potential inconsistencies, but it does not verify. If you ask it "Is this correct?" it may answer with confidence even when it is wrong. The failure mode is "outsourced diligence": the team stops doing the hard work because the tool seems to have done it.

**Misconception 3: "The model is authoritative."** This appears when teams equate polished language with correct reasoning. Consulting culture rewards crisp articulation. The model is very good at crisp articulation. The failure mode is "persuasive error": the document persuades readers to believe something that has not been verified, because it reads like something that would have been verified.

**Misconception 4: "The model is neutral."** This appears when teams forget that the model's output is shaped by the phrasing of the prompt and by the style patterns it has learned. A model can implicitly push toward conventional answers: "best practices," "standard frameworks," and "common recommendations." The failure mode is "default thinking": the organization mistakes generic, conventional output for tailored strategy. The work becomes plausible but unoriginal, and worse, it may not fit the client's reality.

**Misconception 5: "The model can be treated like a junior consultant."** This appears when teams assign the model tasks that require contextual judgment, such as prioritizing initiatives, estimating feasibility, or proposing a recommendation. The failure mode is "responsibility leakage":

because the model produced a recommendation, no one feels fully accountable for it. In a client setting, this is reputationally dangerous. In an internal corporate setting, it can be operationally dangerous.

These misconceptions are amplified by a particular psychological trap: in consulting, the first draft often becomes the anchor. Humans adjust around it. If the first draft contains a hidden assumption, the rest of the work may unconsciously accept that assumption. If the model generates a plausible structure with an implied conclusion, the team may start gathering evidence in support of it rather than testing alternatives. The tool does not merely save time; it shapes the trajectory of thinking. That is why governance must begin with the mental model.

A governance-first team therefore treats the model's fluency as a risk factor. The more polished the output looks, the more careful the team must be. This is counterintuitive for many professionals. Typically, polish signals readiness. With AI, polish signals only that the model can write, not that the content is correct. Level 1 discipline must invert the instinct: a polished draft is still a draft, and may be wrong in subtle ways precisely because it is polished.

One practical way to counter the dangerous misconception is to make the model's ignorance explicit in the output format. That is why Level 1 requires a section called `open_questions`. The presence of open questions forces the reader to see uncertainty. Another is to force the model to list `facts_provided`. This section acts as an "evidence boundary": it shows what the model is allowed to rely on. If a statement appears in the draft output that is not supported by the facts provided, it is immediately suspect.

In other words, the governance wrapper is not bureaucratic overhead. It is the mechanism that prevents the mind from attributing authority to a fluent machine.

### 1.2.3 Definition of a good Level 1 output

Once the mental model is clear, the next step is to define quality. A good Level 1 output is not "a great memo." It is *a draft memo that is epistemically disciplined*: crisp, structured, and explicit about what is known, what is assumed, and what remains unknown. It should be easy to review, easy to verify, and easy to revise.

A useful definition is: *a crisp draft plus explicit unknowns plus traceable inputs plus a clear review path.* Each component matters.

**Crisp draft.** The output should be readable and professionally structured. It should reflect common consulting conventions: clear headings, short paragraphs, bullet points where appropriate, and a coherent narrative arc. This is the part that most teams naturally focus on because it is visible. But crispness is not the main governance objective; it is simply the productivity benefit.

**Explicit unknowns.** A Level 1 output must surface what is missing. It must not conceal uncertainty behind generic filler. Unknowns should be written as questions, and ideally categorized by who must answer them (client, finance, operations, legal, data team). In a consulting setting, explicit unknowns are valuable because they become the agenda for the next client interaction and the internal workplan. In corporate strategy, explicit unknowns become the basis for a verification

plan: what research is needed, what data must be pulled, what stakeholders must be consulted.

**Traceable inputs.** The output must reflect only what was provided. A reader should be able to see the boundary of evidence. This is why `facts_provided` is required. The model should restate the facts succinctly and accurately. If a fact is ambiguous in the input, it should be flagged. Traceability is the beginning of defensibility. In mature environments, traceability later becomes audit artifacts. At Level 1, it is simply a discipline: show your inputs.

**Review path.** A Level 1 output must make review straightforward. That means it should be easy to identify what needs checking. The output should be labeled `"Not verified"` and should include `questions_to_verify`. The review path is the bridge between drafting and professional responsibility. Without it, AI adoption becomes a one-way street: the model produces a draft and it gets sent. With it, AI adoption becomes a controlled loop: draft, review, verify, revise, approve.

This definition implies a practical test: if a reviewer receives the output, can they quickly answer three questions?

1. What do we actually know, and where did it come from?
2. What are we assuming, and who owns those assumptions?
3. What must we verify before this can be sent or used to make a decision?

If the output makes these questions easy to answer, it is a good Level 1 output. If it makes them hard, it is a risky output, regardless of how well it reads.

The definition of "good" also includes what the output must *not* do. A Level 1 output must not:

- Introduce invented numbers, benchmarks, or market facts.
- Imply that verification has occurred.
- Present recommendations as conclusions when they are actually assumptions.
- Hide uncertainty behind vague language that sounds confident.

A good Level 1 output is honest about uncertainty. It can still be persuasive, but it must be responsibly persuasive: persuasion grounded in what is known and explicit about what is unknown.

Finally, a good Level 1 output respects scope. Level 1 is not where the model "solves" the problem. It is where the model helps you articulate the problem and communicate within a disciplined frame. In practice, the most valuable Level 1 outputs are often those that make the next human step clearer: the next meeting, the next analysis, the next data pull, the next decision.

The following artifact box formalizes that definition into a minimum schema. This schema is not optional at Level 1. It is the basic mechanism that makes chatbot drafting governable in consulting environments, where a polished draft can otherwise travel farther and faster than it deserves.

---

**Artifact (Save This)**

**Non-negotiable rule. Facts are not assumptions.** Facts must be provided or verified; assumptions must be stated, owned, and testable.

**Required fields for Level 1 outputs (minimum).**

a. facts_provided (only what you supplied)
b. assumptions (explicit, owned)
c. open_questions (what must be verified)
d. draft_output (the memo/email/outline itself)
e. verification_status: `"Not verified"`

**How to use this schema (practical discipline).**

1. If a statement in `draft_output` is not supported by `facts_provided`, it must be moved to `assumptions` or deleted.
2. If a decision depends on missing information, it must be written in `open_questions` rather than guessed.
3. If the draft will be shared beyond the immediate team, add `questions_to_verify` and name the likely owner of each verification step.
4. Preserve v0 (model draft) and track edits in v1 (human revision) so the team can demonstrate what changed and why.

---

This mental model is intentionally conservative. It does not deny the power of modern language models. It simply refuses to confuse language power with epistemic authority. In consulting and corporate strategy, the cost of that confusion is not a technical error; it is a professional error. Level 1 is where teams build the habit of separating fluency from truth. Once that habit is stable, higher levels of capability can be introduced with proportionate governance. Without that habit, higher levels become a faster way to be wrong.

## 1.3 What Level 1 can do, and what it cannot do

Level 1 is the most deceptively simple stage of the maturity ladder. Because it is "just drafting," teams often assume governance is optional. In practice, Level 1 is where adoption habits form, where cultural norms about AI are created, and where the most common professional failure mode appears: a fluent draft is mistaken for a verified deliverable. The purpose of this section is therefore not to list generic capabilities, but to define a safe scope of use that creates real value without smuggling in hidden risk. In consulting and corporate strategy, the line between drafting and analysis is thin. A well-written paragraph can imply conclusions, certainty, and evidence even when none exists. Level 1 must keep that line visible.

A good Level 1 policy is not "use AI for whatever helps." It is "use AI for these drafting tasks, under these constraints, with these outputs labeled as drafts." That policy allows teams

to capture speed and clarity while preserving professional responsibility. The guiding principle is straightforward: at Level 1, the model may improve *expression*, but it may not become the source of *substance*. Put differently, Level 1 is allowed to make the work more readable, not more true.

### 1.3.1 CAN: drafting tasks that benefit from structure

At Level 1, the highest-return use cases are those where the team already has the substantive content, but needs help turning that content into a coherent artifact. The model is valuable as a formatter, organizer, and tone controller. It is particularly useful when work happens under time pressure, when multiple stakeholders must be aligned, and when consistency matters more than novelty. Below are the canonical "CAN" use cases, with the governance posture that makes them safe.

**Convert rough notes into a clear email, memo, or meeting summary.** This is the most common and most defensible Level 1 use case: you provide the model with redacted notes, and you ask it to produce a structured draft. In consulting, this includes post-meeting follow-ups that capture decisions and next steps, internal updates to partners, and client-facing summaries that confirm what was discussed. In corporate strategy, it includes summaries of stakeholder interviews, quick alignment notes after steering committee meetings, and short memos that frame decisions.

The value is not merely polish. The value is that structure reduces ambiguity. A model can impose a disciplined template: context, what was decided, what remains open, actions with owners, and dates. In doing so, it can reduce the "memory drift" that happens between meetings and documentation. However, the control requirement is strict: the model must use only the facts you provide. If the notes are incomplete, the output must surface open questions rather than fill gaps. A safe Level 1 meeting summary often contains fewer assertions than a typical human summary, because it refuses to invent certainty.

A practical rule is: if you did not hear it, do not let the model write it as if you did. If you are unsure, record it as an open question. That rule protects the team from sending a "confirmation email" that accidentally asserts commitments no one agreed to.

**Produce multiple tones while keeping content constant.** Consulting and strategy work is full of "same content, different audience." A message to a CEO must be concise and decisive; a message to a functional leader may need more context; a message to a skeptical stakeholder may need diplomatic framing. Level 1 is well suited to tone and voice variations because it changes *how* something is said without changing *what* is being said.

The risk here is subtle: tone shifts can introduce content drift. A more "executive" version may compress nuance and drop caveats. A more "persuasive" version may add implied certainty. The control is therefore to treat tone rewriting as a transformation with invariants. The invariant is the factual content boundary: the model may not add new facts, new numbers, or new conclusions. The best practice is to provide the original draft and instruct: "Rewrite into three tones. Do not change meaning. Do not add facts. Preserve all caveats." Reviewers should then compare variants to ensure that disclaimers and open questions remain intact.

Used correctly, this capability improves stakeholder management. It allows teams to speak with clarity without becoming sloppy about uncertainty. It also reduces the temptation to "over-edit" manually, which can introduce inconsistencies across versions.

**Create consistent sections in standard artifacts.** Consulting deliverables often fail not because the underlying thinking is wrong, but because the structure is inconsistent. A memo with no clear objective becomes a narrative dump. A meeting record without actions becomes a forgotten conversation. A draft that lacks risks and dependencies creates false confidence. Level 1 chatbots are particularly useful as "structure enforcers": they can take a set of inputs and reliably produce the same section headings every time.

Common section sets include:

- **Context and objective:** Why are we writing this and what decision is needed?
- **Current state (facts provided):** What is known based on provided inputs?
- **Assumptions and constraints:** What are we taking as true for this draft?
- **Options (if requested):** What alternatives exist, explicitly labeled as unverified.
- **Risks and dependencies:** What could break, and what must be true?
- **Next steps:** What happens next, with owners and timing.

The governance benefit of consistent sections is that they make review systematic. A reviewer knows where to look for assumptions. A manager knows where to find open questions. When teams adopt standard sections, they reduce the chance that critical uncertainties remain hidden in prose. In a governance-first approach, structure is not aesthetics; it is risk management.

**Generate checklists, interview guides, and workshop agendas from your inputs.** This is the other major Level 1 value: producing scaffolds for human work. A consultant preparing for stakeholder interviews can ask the model to generate a question list tailored to the stated objectives. A strategy team designing a workshop can ask for an agenda that fits a given timebox and desired outputs. A transformation team can ask for a checklist of common risks to validate in an implementation plan.

These outputs are safe when treated as *prompts for human inquiry*, not as conclusions. An interview guide is inherently a starting point. A checklist is a reminder system. A workshop agenda is a facilitation structure. None of these require the model to assert facts about the world. They require the model to produce plausible, structured questions and steps. In governance terms, they are "low-claim artifacts": they do not introduce factual assertions that could be wrong; they introduce questions that help humans find the truth.

However, even here, discipline matters. A checklist can still smuggle assumptions, especially if it implies that certain risks are present. The safe posture is to label items as "considerations" and to keep them generic unless the user provided specific facts that justify specificity. If the team wants a checklist "for a cost transformation in a unionized plant," and that is a fact provided, the checklist can include labor relations considerations. If unionization was not provided, the checklist should not assume it.

In summary, Level 1 "CAN" tasks are those where the model improves structure, clarity, and

completeness of *expression*, while leaving the burden of truth entirely with humans.

### 1.3.2 CAN'T: things you must not delegate to Level 1

The "CAN'T" list is where governance becomes real. These are the tasks that are tempting precisely because they appear adjacent to drafting. A model that can write a memo about market entry can easily be asked, "What is the market size?" A model that can draft a board pre-read can easily be asked, "What are the industry benchmarks?" A model that can create a risk section can easily be asked, "What are the regulatory requirements?" At Level 1, those questions are out of scope, because they require verified external truth, not structured language.

**Fact-finding, benchmarking, market data, or citations unless you provide them.** The most common Level 1 failure is invented data presented as if it were researched. In consulting culture, numbers carry authority. A single "typical SG&A reduction is 10–15%" statement can influence a CFO's expectations and lock in a narrative that the team cannot support. Level 1 must therefore prohibit "research-like" outputs unless the underlying facts are explicitly provided by the user and treated as inputs.

This does not mean teams cannot discuss market data. It means the model cannot be the origin of market data. If the team has a market report, or internal data, or a vetted benchmark set, they can provide summarized facts to the model and ask it to incorporate them into a draft. But the model cannot be asked to generate the facts. And if citations are required, they must be human-verified and inserted by the human. The safest rule is: *if you cannot point to the source, the model cannot claim it.*

**Legal, tax, accounting, or regulatory conclusions unless reviewed by experts.** Consulting and strategy often intersect with regulated domains. A restructuring memo may touch employment law. A market entry plan may touch licensing. A cost transformation may touch accounting treatment. Level 1 chatbots can help draft communications and structure issues, but they must not be treated as providers of professional advice. Even when a model's answer sounds plausible, it is not accountable to professional standards, and it may be wrong in jurisdiction-specific ways.

The correct Level 1 posture is to treat these topics as *interfaces* to expert review. The model can help draft questions to counsel, draft an email requesting review, or outline what information might be needed. It can also help rephrase expert guidance that has already been provided, as long as the meaning is preserved and the guidance is quoted or summarized accurately. But it must not invent compliance requirements or interpret rules as if it were licensed to do so. When such content is unavoidable, it must be explicitly labeled as unverified and routed to the appropriate expert for sign-off.

**Final recommendations presented as truth without explicit assumptions.** A core professional obligation in consulting and strategy is to be clear about the basis of recommendations: what evidence supports them, what assumptions they rely on, and what uncertainties could change the conclusion. Level 1 chatbots are very good at writing recommendations with confident language.

This is precisely why they should not be used to create final recommendations at this stage. The risk is not only that the recommendation may be wrong, but that the recommendation may appear justified even when the underlying reasoning and evidence have not been developed.

At Level 1, recommendations can be drafted only under strict constraints: they must be framed as *draft hypotheses* or *options to consider*, accompanied by explicit assumptions and open questions. A safe Level 1 output might say: "Draft recommendation (not verified): pursue Option A assuming X and Y; verify Z before finalizing." What is prohibited is a Level 1 output that reads like a concluded advisory position: "We recommend Option A because it is best practice," without evidence and without assumptions.

This constraint protects the team and the client. It prevents premature closure. It also preserves the integrity of the consulting process: analysis precedes recommendation, not the other way around.

**Confidential data handling outside approved policies and environments.** Level 1 is often adopted through convenience: a public interface, a personal account, or an unvetted tool. This creates a risk that is not about content quality, but about data governance. Consulting engagements and corporate strategy work often involve sensitive information: financial performance, pricing, employee matters, acquisition plans, and proprietary operational data. Even when the model provider has strong security, organizational policy may require specific environments, access controls, and logging.

Therefore, Level 1 "CAN'T" includes any behavior that bypasses approved controls. If the environment is not approved for confidential inputs, do not paste confidential inputs. If the tool does not provide auditability, do not use it for deliverables that require traceability. If the organization requires redaction or anonymization, comply. In a governance-first model, the convenience of drafting is never a justification for violating confidentiality obligations.

A practical principle is minimum-necessary input. Most drafting tasks do not require raw confidential detail. Replace names with placeholders. Remove identifiers. Use ranges instead of exact numbers when feasible. Keep the mapping offline. If the draft truly requires sensitive facts, the work should be done only in an approved environment with appropriate controls.

**A note on "gray zone" tasks.** Many real tasks sit between drafting and analysis. For example, generating "a list of potential risks in this market entry" may look like analysis, but it can be framed as drafting if handled correctly: as a list of considerations to validate, not as claims about reality. Similarly, generating "a hypothesis tree" can be Level 1 if it is treated as a template and explicitly marked as unverified. The difference is not the surface form of the output; it is the epistemic posture. If the output is presented as truth, it is out of scope. If the output is presented as a structured draft to be verified, it can be within scope.

The discipline is therefore to continually ask: *What claims are being made here?* If the model is making factual claims about the world, that is not Level 1 unless those facts were provided. If the model is proposing structure, questions, or ways of organizing known facts, that is Level 1.

**Why the boundaries matter.** These boundaries protect not only against error, but against cultural drift. If a team allows small invented facts to slip in occasionally, the practice will expand.

Over time, the model becomes a quiet source of "knowledge," and the organization loses track of what has been verified. The ladder exists to prevent that drift. Level 1 is where the organization learns to say: "This is drafting, not research," and to enforce that distinction consistently.

In summary, Level 1 is powerful when it is treated as a governed drafting capability. It can make communications clearer, faster, and more structured. It can help teams surface unknowns and organize work. But it cannot be allowed to become a substitute for verification, expertise, or accountable judgment. The "CAN" list captures the value. The "CAN'T" list preserves professional integrity. Together, they define a safe operating envelope that enables adoption without eroding trust.

## 1.4 The Level 1 governance wrapper (minimum controls)

Level 1 governance is not an enterprise bureaucracy and it is not a "later" concern. It is the minimum set of controls that makes chatbot drafting usable in professional consulting and corporate strategy work without degrading trust. The wrapper exists for a simple reason: a Level 1 chatbot is capable of producing text that looks client-ready even when it is unsupported, unverified, or subtly wrong. In other words, the tool is most dangerous when it is most helpful. The governance wrapper is the mechanism that keeps the organization from confusing *fluency* with *reliability* and from confusing *speed* with *quality*.

A common adoption pattern is that teams start informally: someone uses a chatbot to rewrite an email, then to summarize a call, then to draft a one-page memo, and soon the model is producing content that shapes stakeholder expectations. Nothing catastrophic happens at first, which creates a false sense of safety. Over time, however, small errors compound. A phrase that implies a commitment becomes an assumed plan. A guessed number becomes a "benchmark." An assumption hidden in prose becomes a decision premise. A confidential detail pasted in a hurry becomes a policy breach. Level 1 governance exists to prevent these errors from becoming normalized. It is not designed to eliminate every risk. It is designed to eliminate the *avoidable* risks that arise from undisciplined use.

There is also an important cultural purpose. Consulting and corporate strategy are professions of disciplined judgment. If AI adoption at Level 1 teaches teams that "outputs are disposable and unowned," the profession weakens. If Level 1 adoption teaches teams that "outputs are drafts with explicit uncertainty and traceable inputs," the profession strengthens. The wrapper therefore functions as a training system. It makes the right habits automatic: facts are listed, assumptions are owned, open questions are explicit, and a human review gate is expected.

The minimum control set is intentionally small and enforceable. If the wrapper requires too much overhead, teams will bypass it. If it is too vague, teams will interpret it inconsistently. The goal is a minimal discipline that can be applied on real engagements at real speed.

---

**Risk & Control Notes**

**Capability.** Drafting acceleration for consulting artifacts

**Primary risks.** Persuasive error; hidden assumptions; leakage of client-sensitive info; false precision; loss of authorship accountability

**Minimum controls.** Use a standard input/output schema; redact sensitive inputs; label outputs as drafts; maintain a prompt log; require human review and sign-off; keep version history

---

The above summary can be read as a compact operating model for Level 1. Each element of the minimum controls deserves operational clarification.

**1) Use a standard input/output schema.** The schema is the most important Level 1 control because it directly targets the primary failure mode: the blending of facts and assumptions into persuasive prose. A schema forces the model to separate what it was given from what it is inferring. In governance terms, it creates an "epistemic boundary." In practical terms, it gives reviewers a fast way to scan the draft for risk: if the `facts_provided` section is thin but the `draft_output` is confident, something is wrong.

The schema also helps in a less obvious way: it makes the model a tool for surfacing uncertainty rather than hiding it. The requirement to list `open_questions` creates a natural feedback loop. Every time the model produces an open question, the team can assign it, verify it, and update the draft. Over time, teams learn to treat drafting as an iterative process that moves from unknowns to knowns.

A good Level 1 schema is short and strict. It should include at minimum:

- `facts_provided`: what the user supplied (not what the model thinks is true).
- `assumptions`: explicit premises the draft relies on.
- `open_questions`: missing information that must be verified or obtained.
- `draft_output`: the artifact itself.
- `verification_status`: "Not verified".

The schema becomes the backbone of defensibility. If later someone asks, "Why did we say this?" the team can point to inputs, assumptions, and review notes instead of relying on memory.

**2) Redact sensitive inputs (minimum-necessary).** Redaction is often treated as an IT policy issue. In practice, it is a professional discipline issue. Consultants and strategy teams routinely handle information that is sensitive not only because of legal confidentiality, but because of competitive impact: pricing strategy, margin structure, acquisition intentions, labor relations, operational vulnerabilities, and internal performance metrics. Even if the tool environment is secure, organizations may have rules about what can be shared, where it can be shared, and how it must be logged.

The Level 1 principle is minimum-necessary input: provide the model only what it needs to perform the drafting task. Most drafting tasks require far less detail than teams assume. A follow-up

email does not require a client name; "Client A" is enough. A decision memo does not require the exact revenue number; "low single-digit growth" may be enough at the drafting stage. A meeting summary rarely needs full names; roles are sufficient ("Head of Ops," "CFO," "PMO Lead").

Redaction is not about paranoia. It is about reducing exposure without reducing usefulness. The less sensitive input you provide, the lower the consequence of tool misuse or policy mismatch. Redaction also supports reuse: a redacted prompt can be saved as a template and shared internally without risking leakage.

**3) Label outputs as drafts.** This is the simplest control and the most culturally important. If outputs are not explicitly labeled as drafts, they will be treated as final. People forward emails. They paste paragraphs into decks. They move quickly. In that environment, the phrase "Not verified" is not a disclaimer; it is a traffic signal. It tells the organization: stop, review, verify.

This label must not be optional. It should be present in every Level 1 output. It should also be reinforced by the workflow: outputs are stored as v0 drafts, reviewed by a human, revised into v1, and only then used externally. The label protects recipients from misinterpreting the content, and it protects authors from accidentally over-claiming.

**4) Maintain a prompt log.** A prompt log is often misunderstood as surveillance. In governance-first practice, it is recordkeeping. Consulting work relies on evidence of process. When a deliverable is questioned, the ability to reconstruct how it was produced matters. At Level 1, the prompt log can be minimal: the prompt text, the redacted input snapshot, and the output. This creates a lightweight audit trail.

The prompt log also has a learning function. Over time, teams identify which prompt patterns are safe and effective. They can refine templates, standardize language, and reduce variability. This is how Level 1 usage becomes internal capability rather than individual improvisation.

A practical point: the prompt log should support redaction and should avoid storing sensitive information unnecessarily. The point is to be able to demonstrate the framing and constraints used, not to create a repository of confidential content.

**5) Require human review and sign-off.** Human review is the control that enforces accountability. Without it, Level 1 becomes a pipeline from the model to the outside world. With it, Level 1 becomes a drafting step inside a human-owned process. The review does not need to be heavy for every artifact, but it must exist in proportion to risk.

In practice, organizations can define tiers:

- **Low-risk internal drafts** (e.g., internal brainstorming email): self-review by the author.
- **Moderate-risk stakeholder communications** (e.g., client follow-up, exec update): peer review or manager review.
- **High-risk communications** (e.g., board materials, public statements, sensitive negotiations): senior review and explicit sign-off.

The key is that the organization knows which tier applies and can show evidence of review. In consulting, "sign-off" can be as simple as a note in a tracked document or a comment in a ticketing system. The form is less important than the reality: someone is responsible.

**6) Keep version history.** Versioning is the bridge between AI drafting and professional authorship. If you keep v0 (model output) and v1 (human edits), you can demonstrate what the model contributed and what the human corrected or verified. This protects the team in two ways. First, it creates transparency: the human is not pretending the model did not exist; the team is showing that it was used responsibly. Second, it improves quality: when teams can see patterns in what they consistently correct, they can refine their prompts and templates.

Version history also supports controlled reuse. A v0 that consistently requires heavy editing is not a good template. A v0 that reliably produces a safe structure may be promoted into a standard prompt pattern.

With those controls clarified, the next step is to define the concrete unit of work that operationalizes them: the Draft Bundle.

### 1.4.1 The "Draft Bundle" (what you must save)

The Draft Bundle is the minimum package of artifacts that makes Level 1 work defensible. It is not meant to slow teams down. It is meant to preserve evidence of responsible use. In consulting and corporate strategy, deliverables live longer than memory. Months later, a client may ask why a statement was made. An internal audit may ask how a message was produced. A partner may need to reconstruct the logic behind a recommendation. Without artifacts, the answer becomes guesswork. With a Draft Bundle, the answer becomes process evidence.

The Draft Bundle also normalizes a critical professional habit: treating drafts as part of a traceable workflow rather than as ephemeral text. This habit is what later enables Level 3 and Level 5 governance. If teams cannot save a Draft Bundle at Level 1, they are not ready for agentic workflows.

---

**Artifact (Save This)**

**Minimum deliverable standard (Draft Bundle).**

1. Input snapshot (redacted): what facts you provided to the model.
2. Prompt text: the exact instruction used.
3. Output JSON: facts/assumptions/open_questions/draft_output + *Verification status: Not verified. Human review required.*
4. Reviewer notes: what changed, what was verified, what remains open.
5. Version history: v0 (model), v1 (human edits), v2 (final draft).

---

Each element of the Draft Bundle is small, but together they create defensibility.

**Input snapshot (redacted).** This is the evidence boundary. It answers the question: what did we tell the model? Without this, you cannot evaluate whether the model invented content. The snapshot can be short: bullet facts, not full transcripts. It should also be redacted to remove sensitive identifiers. The goal is not to preserve everything; the goal is to preserve what matters to interpret the output.

**Prompt text.** This is the instruction boundary. It answers the question: what constraints did we impose? A prompt that says "Write a persuasive memo" is different from a prompt that says "Use only the facts provided; do not invent; list assumptions and open questions." If the organization is serious about governance, it must be able to show that it used constraints consistently.

**Output JSON.** A structured output is not merely a formatting preference. It is a governance artifact. It forces separation of facts and assumptions, and it makes verification status explicit. It also makes the output machine-readable for later aggregation, search, and audit. Even at Level 1, a consistent JSON schema enables later governance improvements: sampling, QA, and pattern detection (e.g., identifying prompts that often produce missing assumptions).

**Reviewer notes.** This is the accountability evidence. It answers: what did the human change, and what did they verify? Reviewer notes should be short but explicit. They do not need to be essays. A few bullet points are sufficient:

- Corrected an overstated claim about timeline.
- Removed an implied benchmark that was not provided.
- Verified revenue range against internal finance dashboard.
- Left open questions regarding data readiness for next meeting.

These notes demonstrate professional discipline. They also train teams to review AI drafts critically.

**Version history.** This is the provenance trail. It answers: what evolved from draft to final? In many organizations, versioning already exists through document tools. The Level 1 requirement is simply to preserve the model draft as v0. This prevents "memory laundering," where the final document no longer shows what was AI-generated and therefore cannot be audited or learned from. Version history also supports coaching: leaders can review how juniors are using AI and where their judgment needs strengthening.

The Draft Bundle can be implemented lightly. It can be a folder with five files. It can be a single document with sections. The format is flexible. The requirement is not.

### 1.4.2   Redaction and minimum-necessary input

Redaction is the most immediate and practical control in consulting environments because it protects confidentiality regardless of model behavior. Even a perfectly accurate model can create risk if the user inputs sensitive information into an inappropriate channel. The Level 1 governance wrapper therefore requires redaction not as a "best practice," but as a default.

Minimum-necessary input means you should provide only what the model needs to produce the draft. That sounds obvious, but it is frequently violated because teams copy and paste entire source materials to "save time." The governance-first posture asks you to slow down for one moment: what is the minimum information required to draft this artifact? Often it is a short set of bullets:

- Audience and purpose (who is this for and what decision is needed).
- A few factual points (what happened, what was agreed, what constraints exist).
- Tone and format requirements (email, memo, summary; executive vs detailed).

Most other details can be added later by humans in approved environments.

Redaction is also a way to standardize and productize Level 1 prompts. If your prompts contain confidential specifics, they cannot be shared or reused. If your prompts are redacted and templated, they become internal assets that improve team quality.

There is an additional governance benefit: redaction forces you to clarify the underlying logic. When you replace "Acme Corp acquisition of BetaCo" with "Client A potential acquisition of Target B," you are forced to think about what is structurally relevant versus what is merely contextual. This often improves the quality of the draft.

The following checklist provides a minimum set of redaction rules that are simple enough to follow under time pressure.

---

**Checklist**

**Level 1 redaction checklist (minimum).**

1. Replace client, product, and counterparty names with placeholders (Client A, Segment B).
2. Remove unique identifiers (emails, phone numbers, account IDs, contract IDs).
3. Do not paste raw financial statements, HR records, or source code unless approved.
4. Use ranges or normalized values (e.g., "mid-teens margin") where feasible.
5. Keep a local mapping file (offline) if you must re-identify later.

---

To make the checklist operational, teams should adopt a few concrete practices.

**Use role labels instead of names.** Replace "Maria Gonzalez" with "Head of Procurement." Replace "John Smith" with "CFO." In most drafting tasks, the role is what matters. This practice also improves clarity, because roles are more meaningful to readers than names, especially in early drafts.

**Abstract deal terms.** Replace exact prices or terms with ranges or placeholders: "purchase price in the low billions," "term sheet under negotiation," "target close in Q3 (subject to approvals)." The model does not need exact numbers to draft a structure. If exact terms must be included, they should be inserted by humans later in the controlled document.

**Remove embedded documents.** Do not paste raw statements, contracts, or internal reports unless policy explicitly permits and the environment is approved. Instead, summarize the relevant facts in bullet form. This has a quality advantage: it forces the human to decide what is relevant. It also has a governance advantage: it limits data exposure.

**Maintain offline mapping.** If placeholders are used, maintain a local mapping file that is not shared with the model. This ensures you can later re-identify for final drafting while preserving confidentiality during AI-assisted drafting.

**Treat redaction as part of the workflow, not an optional step.** The easiest way to ensure compliance is to make redaction part of the template: every Level 1 prompt begins with "Replace names with placeholders." Over time, this becomes automatic.

**Escalate when in doubt.** When a user is unsure whether information is sensitive, the

governance-first approach is to escalate rather than guess. That may mean using a more controlled internal environment, or it may mean excluding the information from the AI prompt and inserting it manually later.

The Level 1 governance wrapper is, ultimately, a professional discipline system. It makes AI drafting safer by forcing transparency about knowledge boundaries, by protecting confidentiality through minimum-necessary input, and by preserving accountability through review and versioning. It also sets the stage for later levels. When teams can reliably produce Draft Bundles, maintain prompt logs, and enforce redaction, they have demonstrated the foundational control maturity needed for Level 2 reasoning scaffolds and, eventually, Level 3 workflow orchestration.

If Level 1 governance feels "heavy," that is usually a sign that the organization is trying to use Level 1 for tasks that belong to higher levels or to human expertise. Used correctly, Level 1 governance is light: a schema, a label, a log, a review, and a version trail. These are not new burdens. They are the professional habits that consulting and strategy already rely on. Level 1 simply makes those habits explicit in a world where language is easy and therefore mistakes can travel fast.

## 1.5   Core Level 1 workflow patterns (A–D)

Level 1 becomes genuinely useful when it is not treated as improvisation, but as a repeatable workflow. The goal is not to produce "good prompts." The goal is to produce consistent, auditable drafting behaviors that teams can rely on under real delivery pressure. In consulting and corporate strategy work, the same drafting moments occur repeatedly: after a meeting, before a steering committee, when translating analysis into a narrative, and when shaping an executive storyline. Pattern libraries exist in every strong team, even if they are informal. Level 1 simply makes them explicit and governable.

Each workflow pattern below is designed to meet four requirements simultaneously:

1. **Speed:** it must work under time pressure and with imperfect notes.
2. **Structure:** it must force a consistent output shape that is easy to review.
3. **Epistemic discipline:** it must prevent the model from inventing facts and must surface uncertainty.
4. **Governance trace:** it must naturally produce artifacts suitable for the Draft Bundle (inputs, prompt, output, reviewer notes, versions).

The patterns are written to be used by consulting teams and corporate strategy functions with minimal AI background. They do not assume tool integrations, agents, or autonomous research. They assume a single chatbot session operating under a strict "no new facts" rule, producing drafts labeled as not verified. They also assume that humans, not the model, own the final message.

### 1.5.1 Pattern A: Notes-to-email (stakeholder follow-up)

**When to use.** Pattern A is used immediately after a meeting, call, or workshop when the team needs to send a follow-up email to confirm what was discussed, document decisions, and align on next steps. This is one of the highest-return Level 1 patterns because it reduces the time cost of writing and reduces the risk of misalignment. It is also one of the highest-risk patterns if used carelessly, because a follow-up email can accidentally create commitments. A governance-first follow-up email must therefore prioritize accuracy, explicit uncertainty, and careful wording.

**Inputs.** The inputs are short bullet notes. The notes should include only what the team is confident was said or decided, plus any clearly stated action items. If the notes are uncertain, they should be marked as uncertain. Redaction should be applied before sharing with the model: replace names with roles, replace client names with placeholders, and avoid unique identifiers.

A minimal input set looks like:

- Audience: who will receive the email (e.g., CFO + Head of Ops).
- Purpose: confirm discussion + decisions + actions.
- Facts observed: what was discussed, what was agreed.
- Actions: who will do what by when (if known).
- Open items: what was unclear, what requires follow-up.
- Tone: neutral executive, diplomatic, or direct.

**Output.** The output is a structured email draft plus a separate list of open questions that must be verified before sending. In Level 1 governance, it is good practice to keep open questions out of the email body if they are not appropriate to send externally, but to list them for internal review. Alternatively, some open questions may be included in the email as "confirmations" if they are intended to be resolved with the recipient.

**Core logic.** The key risk in stakeholder follow-up is "implicit commitment." The model may rephrase a vague statement into a firm-sounding commitment. The control is to constrain the language: use confirmation language rather than promise language; use "as discussed" rather than "we will"; and separate what was agreed from what is proposed. If a timeline is not agreed, it should not appear as a date. If a deliverable is not promised, it should not be written as if it is.

**Step-by-step workflow.**

1. **Redact and normalize notes.** Convert names to roles; remove identifiers; mark uncertain items with "(uncertain)" or "(to confirm)."
2. **Force the model into a strict schema.** Require `facts_provided`, `assumptions`, `open_questions`, `draft_output`, `verification_status`.
3. **Specify email constraints.** Ask for a subject line, short paragraphs, and explicit next steps.
4. **Review for commitment language.** Humans must scan for words like "will," "commit," "guarantee," "by Friday," and replace with conditional language where needed.
5. **Send only after verification.** If the email includes factual claims (numbers, dates, decisions), confirm them against notes or with meeting participants.

**Common failure modes and how to prevent them.**

- **Failure: the email implies agreement that did not exist.** Prevention: include a line such as "Please reply if any of the below does not reflect your understanding."
- **Failure: the email introduces a timeline.** Prevention: instruct "Do not introduce dates unless provided. If timing is unclear, state 'timing to be confirmed.'"
- **Failure: the email sounds too confident.** Prevention: use "as discussed," "initial view," "draft," and keep recommendations out unless explicitly requested.

**What good looks like.** A good Level 1 follow-up email is short, accurate, and explicit about next steps. It does not over-claim. It invites correction. It helps the team by converting messy notes into a professional message while keeping uncertainty visible.

### 1.5.2 Pattern B: Meeting transcript-to-summary (decision record)

**When to use.** Pattern B is used when the team needs a decision record: a structured summary of what happened in a meeting, what was decided, what actions were assigned, and what remains open. This is common in consulting engagements with steering committees, and in corporate strategy functions running transformation programs. Unlike a simple summary, a decision record is a governance artifact. It is meant to reduce ambiguity and to provide traceability over time.

**Inputs.** The input may be a transcript, but Level 1 governance recommends avoiding raw transcripts unless the environment is approved and the transcript is already sanitized. In many cases, it is safer to provide the model with a condensed "pseudo-transcript": a sequence of bullet points with speaker roles and key statements. If a real transcript is used, it should be redacted and treated as sensitive.

Inputs should include:

- Meeting context: date, purpose, attendees (roles only).
- Key points raised: major arguments, constraints, concerns.
- Decisions stated: what was explicitly agreed (if any).
- Actions stated: who owns what next step (if any).
- Unresolved items: questions that remain open.

**Output.** The output is a decision record with explicit sections:

- **Decisions:** statements that reflect explicit agreement.
- **Actions:** action items with owners and due dates (or "TBD").
- **Open items:** unresolved questions and dependencies.
- **Risks/notes:** key concerns raised, without inventing.

**Core logic.** The main risk in Pattern B is the model "upgrading" discussion into decision. Humans talk in tentative language: "we should," "maybe," "it sounds like," "I think." A chatbot may convert this into "The team decided." That is unacceptable in a decision record. The control is to define a decision rule: a statement is a decision only if it was explicitly stated as such or if the

chair confirmed it. Everything else belongs in "discussion points" or "open items."

**Step-by-step workflow.**

1. **Pre-process the transcript.** If possible, convert to bullets with time stamps or speaker roles. Remove sensitive details.
2. **Apply the decision rule.** Mark which items are confirmed decisions versus tentative.
3. **Prompt for a decision record format.** Require the model to use your sections and to label uncertain items.
4. **Human validation pass.** The meeting owner reviews the draft against notes and confirms decisions, owners, and due dates.
5. **Publish with correction mechanism.** Include a line: "Please respond with corrections within 24 hours."

**Common failure modes and how to prevent them.**

- **Failure: invented owners or due dates.** Prevention: instruct "If owner or date not provided, write 'Owner: TBD' and 'Due: TBD'."
- **Failure: decisions inferred from tone.** Prevention: instruct "Only label as decision if explicitly stated. Otherwise place under 'Open items'."
- **Failure: omission of key risks.** Prevention: include a "concerns raised" section and ask the model to list them verbatim-style from your bullets.

**What good looks like.** A good decision record is not "complete" in the sense of covering everything. It is complete in the sense of capturing what the organization needs to proceed safely: decisions, actions, and unresolved issues. It is easy to scan, and it creates accountability without overstating certainty.

### 1.5.3 Pattern C: Outline-to-memo (executive narrative)

**When to use.** Pattern C is used when the team has an outline of analysis or a set of key messages, and needs to translate it into an executive narrative: a one-page memo, a pre-read for a steering committee, or a short decision brief. This is a common consulting deliverable and a common corporate strategy artifact. The memo must be readable, logically ordered, and honest about what is known and unknown.

**Inputs.** The input is an outline, not raw analysis. The outline should contain:

- The decision or purpose of the memo.
- Key facts that are verified (provided by the team).
- The central hypothesis or recommendation (if any), clearly labeled as draft.
- Key arguments supporting the hypothesis.
- Key risks, constraints, and dependencies.
- Open questions and data gaps.

In Level 1, the model should not be asked to generate new arguments based on external facts. It

should be asked to convert your outline into a coherent narrative with clear structure.

**Output.** The output should follow a consistent executive structure. A common consulting shape is:

- **Context:** what situation the memo addresses.
- **Key question:** what decision is needed.
- **Draft answer (not verified):** the proposed direction, clearly labeled.
- **Rationale:** the key supporting points based on provided facts.
- **Risks and sensitivities:** what could change the conclusion.
- **Next steps:** what must be verified and what actions follow.

**Core logic.** The main risk in Pattern C is that the model will "tighten" the narrative by smoothing away caveats. Executives like clarity. The model produces clarity. The danger is false certainty. The control is to require explicit caveats as first-class elements, not optional footnotes. If the memo relies on assumptions, they must be listed. If the conclusion depends on unverified data, that dependency must be stated. If the analysis is preliminary, it must be labeled as such.

Another risk is that the model may improve rhetorical flow by adding connective tissue that sounds like evidence. For example, it may write "industry leaders typically do X," even if you did not provide that fact. The control is the "no new facts" rule and a review pass that checks every claim against `facts_provided`.

**Step-by-step workflow.**

1. **Prepare an outline with explicit tags.** Mark items as Fact, Assumption, or Open Question.
2. **Prompt for a memo conversion.** Ask the model to preserve the tags and to present assumptions and open questions explicitly.
3. **Review for invented claims.** Check for phrases like "typically," "industry standard," "on average," "best practice," unless you provided sources.
4. **Strengthen the review path.** Add a "Questions to verify" list that assigns owners (Finance, Ops, Data, Legal).
5. **Finalize through human edits.** Insert verified numbers and sources manually, and update verification status when appropriate.

**What good looks like.** A good executive memo draft reads cleanly while staying honest about uncertainty. It makes the decision crisp, but it does not pretend the evidence base is stronger than it is. It turns an outline into a narrative without turning speculation into fact.

### 1.5.4 Pattern D: Slide storyline (titles only, no numbers invented)

**When to use.** Pattern D is used when the team needs a slide storyline: a sequence of slide titles that tell a coherent "so what" narrative. This is a core consulting skill. A storyline is not a deck. It is the skeleton that later becomes a deck. In Level 1, producing titles only is a powerful control: it captures the narrative arc without inviting the model to fabricate data to fill charts.

**Inputs.** The input is the purpose of the deck and the key messages the team wants to

communicate. It may include:

- Audience and decision context (board update, CEO brief, steering committee).
- Objective (approve direction, align on plan, decide investment).
- Key facts that are verified (high-level, redacted).
- The main storyline thesis (draft, if needed).
- Constraints (time, scope, sensitivity).

**Output.** The output is a list of slide titles (and optionally one-line "speaker notes" per slide) with explicit placeholders where data is missing. It should not include numbers, charts, or claims of benchmarks unless provided. It should not include fake citations. The storyline should be driven by the "so what" principle: every slide title should express a message, not a topic.

For example, "Cost transformation opportunities" is a topic. "We can fund growth by eliminating structural costs in procurement and SG&A" is a message. Pattern D aims for messages, because messages force clarity.

**Core logic.** The main risk in slide drafting is false precision and overreach. Slides invite numbers. The model is tempted to provide them. The governance control is to forbid numbers and to require placeholders: "Insert verified baseline cost by category," "Insert verified timeline," "Insert verified financial impact ranges." This turns the model into a narrative shaper, not a data generator.

Another risk is storyline bias: the model may default to conventional frameworks (market overview, SWOT, 3-year plan) regardless of the actual decision. The control is to force a decision-driven structure: start with the decision needed, then evidence, then options, then recommendation constraints, then risks and next steps.

**Step-by-step workflow.**

1. **Specify the decision and audience.** A storyline without a decision is just slides.
2. **Provide the few verified facts you have.** Keep them high-level and redacted.
3. **Instruct "titles only."** Explicitly ban numbers, benchmarks, and invented data.
4. **Require placeholders for missing evidence.** Every slide that needs data must state what data is required.
5. **Review for hidden claims.** Even without numbers, titles can imply facts. Ensure titles do not assert what is not known.

**What good looks like.** A good storyline is coherent, decision-driven, and honest about what evidence is needed. It accelerates deck building by providing a narrative map, but it does not pretend the supporting data exists. It makes gaps visible early, which is exactly what good consulting teams need.

**How patterns A–D fit together.** These patterns are complementary. A common sequence in real work is:

- Pattern B produces a decision record after a steering meeting.
- Pattern A produces a follow-up email confirming actions to stakeholders.
- Pattern C converts the emerging outline into a short executive memo for the next checkpoint.

- Pattern D creates a storyline for the next steering committee deck.

This sequence shows the real value of Level 1: not isolated drafting tricks, but repeatable communication moves that keep work aligned and auditable.

**A final governance note.** None of these patterns remove the need for professional judgment. They make the human's job easier by giving structure and clarity, but they also increase the importance of review because the output looks polished. Teams should therefore treat every pattern output as v0, preserve it, and apply reviewer notes before anything is distributed. Over time, the patterns become internal standards. That is the correct outcome of Level 1 adoption: a disciplined drafting system that makes professional work faster without making it less defensible.

## 1.6 Mini-cases (Level 1): drafting in real consulting scenarios

The purpose of the Level 1 mini-cases is not to demonstrate that a chatbot can "do consulting." That framing is precisely the mistake this book is trying to prevent. The purpose is to demonstrate that a chatbot can accelerate the production of *draft artifacts* that consulting and corporate strategy professionals already create, while preserving the core disciplines that make those artifacts defensible: clear ownership, explicit assumptions, transparent unknowns, and a review path.

Each mini-case below follows the same structure:

1. A small, explicitly limited fact pattern. These facts are the only inputs the model is allowed to use.
2. The Level 1 drafting task: what artifact we want (memo, email, agenda).
3. The expected output posture: facts separated from assumptions, open questions surfaced, and verification status labeled as not verified.

The reader should view these mini-cases as "templates for practice." In real work, the facts will differ, but the discipline should remain constant. Most importantly, the drafts produced at Level 1 are not final deliverables. They are v0 drafts that must be reviewed and verified by humans before being shared externally or used to support decisions.

### 1.6.1 Case 1: Market entry (one-page decision memo)

**Scenario.** A corporate strategy team at *Client A* (a mid-sized industrial components manufacturer) is evaluating whether to enter *Country B* in Latin America. The CEO has asked for a one-page decision memo to frame the decision for the executive committee. The team has limited time and only a small set of known facts. The memo is needed to align on what must be verified before a recommendation can be made.

**Fact pattern (inputs provided to the model; redacted).**

- Client A manufactures industrial components for heavy equipment OEMs.
- Client A currently sells primarily in North America and has limited presence in Latin America.
- Country B has a growing mining sector and several OEM customers operate there.

- Two entry options are under consideration:

  a. Distributor partnership (local distributor sells and services).
  b. Direct presence (small sales office + service capability within 18 months).

- Constraints:

  – Initial investment must be "modest" (no acquisition in year 1).
  – Leadership is risk-sensitive about compliance and operational execution.
  – The CEO wants an initial decision framework in two weeks.

- What is not yet known:

  – Market size in Country B for Client A's specific product categories.
  – Competitive landscape and pricing dynamics.
  – Regulatory and compliance requirements for service operations.
  – Availability of qualified service technicians and logistics constraints.

**Level 1 drafting task.** Produce a one-page decision memo draft that frames the decision and organizes the work. The memo should not contain invented market numbers or claims about competitors. It should clearly separate facts provided from assumptions, list open questions, and label the output as not verified.

**What a good Level 1 draft accomplishes.** In market entry work, early memos can be dangerous when they overreach. A chatbot can easily generate a confident "recommendation" based on generic frameworks. That is not the goal here. The goal is to produce a crisp decision framing memo that:

- States the decision clearly.
- Summarizes the options without implying a conclusion.
- Lists assumptions explicitly (e.g., what "modest investment" implies).
- Surfaces open questions that must be answered before a recommendation.
- Proposes next steps in a way that supports verification.

**Illustrative Level 1 output shape (not the final answer, but the structure).** A one-page memo draft should typically include:

- Context and objective.
- Decision required and timing.
- Options under consideration (A/B).
- Evaluation criteria (e.g., speed, risk, investment, control).
- Risks and dependencies (explicit).
- Open questions / data required.
- Next steps and owners (if known).

**Governance reminders in this case.** Market entry memos often invite the model to invent market sizing and competitor facts. The Level 1 constraint must explicitly forbid it. The review pass must check the memo for "default consulting claims" such as "the market is expected to grow"

or "competitors typically price at X." If such statements appear without being provided as facts, they must be removed or moved to open questions.

The value of the Level 1 memo is speed and structure: it gives the executive committee a clean frame and gives the team a clear verification agenda. It does not provide the answer. It organizes the path to the answer.

### 1.6.2 Case 2: Cost transformation (workplan + stakeholder email)

**Scenario.** A consulting team is supporting *Client C*, a consumer products company facing margin pressure. Leadership wants a rapid cost transformation diagnostic to identify near-term savings opportunities and a longer-term program roadmap. The consulting team needs two Level 1 artifacts immediately: (i) a draft workplan outline for a four-week diagnostic, and (ii) a cautious kickoff email to internal stakeholders to set expectations without overpromising savings.

  **Fact pattern (inputs provided to the model; redacted).**

- Client C has experienced gross margin compression over the past year.
- Leadership has requested a four-week "rapid diagnostic" focused on:

  - Procurement (direct + indirect).
  - SG&A efficiency.
  - Manufacturing and logistics (high level only; no site visits in first four weeks).

- Constraints:

  - No headcount reduction decisions will be made in the first four weeks (leadership statement).
  - The CFO wants "credible ranges" only after data validation.
  - The organization is sensitive to morale; communications must be cautious.

- Known stakeholders:

  - CFO sponsor.
  - Head of Procurement.
  - Head of Operations.
  - HR partner (for comms alignment).

- What is not yet known:

  - Baseline spend by category and supplier concentration.
  - SG&A breakdown and allocation drivers.
  - Manufacturing cost drivers and major constraints.
  - Current contract terms, renewal timing, and procurement policies.

  **Level 1 drafting task.** Produce:

a. A workplan outline for a four-week rapid diagnostic (phases, workstreams, deliverables, data requests).

b. A kickoff email draft to stakeholders that sets a professional tone, clarifies that outputs are preliminary until data validation, and requests initial data/support.

Both artifacts must be drafts labeled not verified and must not include invented savings targets or benchmarks.

**What a good Level 1 draft accomplishes.** Cost transformations are high-stakes and politically sensitive. The fastest way to lose credibility is to promise savings before understanding the baseline. A chatbot can help by creating:

- A structured workplan that looks professional and is easy to execute.
- A consistent set of data request categories (without claiming the data exists).
- A kickoff email that communicates seriousness without causing panic.

The workplan should focus on sequencing and governance: weekly milestones, stakeholder touchpoints, and the principle that savings ranges are produced only after validation. The email should be careful about language. It should avoid "we will deliver X% savings." It should instead communicate "we will identify opportunities and validate ranges."

**Governance reminders in this case.** Workplans are often copied forward and become commitments. Therefore:

- If dates are not agreed, use "Week 1 / Week 2" rather than calendar dates.
- If deliverables depend on data availability, state that dependency explicitly.
- For the email, include an invitation to correct misunderstandings and a clear statement that the diagnostic is fact-finding, not decision-making, in the first phase.

A Level 1 chatbot is well suited to this case because both artifacts are primarily about structure and tone, not about substantive savings estimation. Human judgment remains responsible for validating data and shaping actual recommendations.

### 1.6.3 Case 3: Capital allocation (board pre-read outline)

**Scenario.** A corporate strategy and finance team is preparing for a board meeting where the board will discuss capital allocation priorities. The CEO wants a pre-read structure that is rigorous but does not overclaim. The team has partial inputs and needs a draft outline that clearly states what data is required before finalization. The goal at Level 1 is to draft a board-ready structure (not the final content), with placeholders and explicit open questions.

**Fact pattern (inputs provided to the model; redacted).**

- The company has three competing uses of capital:

  a. Growth investment in capacity expansion.
  b. Debt reduction to improve credit metrics.
  c. Shareholder return program (dividend and/or buybacks).

- The board has requested that management present:

  – A clear decision framework and tradeoffs.
  – Scenario sensitivity (base / downside).
  – Risks to liquidity and covenants (high level).

- Constraints and posture:

  - Management wants to avoid committing to a single path before Q2 forecast refresh.
  - The CFO insists that any figures shown must be reconciled to internal financials.
  - There is an expectation of an executive summary plus appendices.

- What is not yet available (for this drafting step):

  - Final forecast numbers and full scenario model outputs.
  - Updated covenant headroom calculations.
  - Project-level returns for capacity expansion options.
  - Share repurchase constraints and timing considerations.

**Level 1 drafting task.** Draft a board pre-read *outline* (headings, subheadings, and a short description of what each section must contain). The outline must:

- Be decision-driven (not a generic "finance update").
- Include explicit placeholders for required data.
- Include an open-questions list that makes verification requirements visible.
- Avoid inventing numbers, returns, or "recommended" policies.

**What a good Level 1 draft accomplishes.** Boards do not need prose; they need clarity. A Level 1 chatbot can accelerate the structuring of a board pre-read by ensuring that:

- The decision question is explicit.
- Tradeoffs are framed in a consistent way across options.
- The document includes the right risk categories (liquidity, covenants, execution).
- Missing data requirements are visible early.

The outline becomes a coordination artifact for the finance and strategy team: it tells them what analysis and validation must be completed before the pre-read is finalized. It also helps prevent a common error: a board document that looks complete but is missing critical sensitivities.

**Governance reminders in this case.** Board materials are high-risk communications. At Level 1:

- Use placeholders aggressively: "Insert verified forecast," "Insert covenant headroom," "Insert project IRR range."
- Include a verification checklist: what must be reconciled, who must sign off (Finance, Treasury, Legal).
- Maintain version control: board materials evolve quickly; preserving v0 and reviewer notes is essential.

The chatbot's role is to structure the pre-read, not to decide capital allocation. Any recommendation language should be framed as "options under consideration" unless management has explicitly decided otherwise.

### 1.6.4   Case 4: Operating model redesign (workshop agenda)

**Scenario.** A consulting team is supporting *Client D* in an operating model redesign for a corporate function (e.g., procurement, finance, or customer operations). The first step is a cross-functional workshop to align stakeholders on current pain points, design principles, and candidate operating model options. The team needs a workshop agenda and an interview guide to use in the week leading up to the workshop. The artifacts must be professional, time-boxed, and designed to elicit factual inputs rather than opinions disguised as facts.

**Fact pattern (inputs provided to the model; redacted).**

- Objective: define target operating model principles and identify 2–3 viable operating model options for evaluation.
- Participants (roles): Functional lead, HR partner, IT partner, Finance partner, Regional representatives, PMO lead.
- Constraints:
  - Workshop duration: 3 hours.
  - The organization has recent change fatigue; facilitation must be constructive.
  - Leadership wants clear next steps and ownership after the workshop.
- Known themes from prior conversations (unverified):
  - Complaints about unclear decision rights.
  - Inconsistent processes across regions.
  - Data and systems fragmentation.
- What is not yet known:
  - Current-state process maps and cycle times.
  - Actual pain point frequency and severity by region.
  - Existing governance forums and decision rights documentation.
  - System constraints and roadmap.

**Level 1 drafting task.** Produce:

a. A 3-hour workshop agenda with timeboxes, objectives per segment, and expected outputs.
b. A semi-structured interview guide (10–15 questions) to run stakeholder interviews prior to the workshop.

Both artifacts must include governance reminders: separate facts from opinions, capture open questions, and avoid presuming solutions.

**What a good Level 1 draft accomplishes.** Operating model redesign work is vulnerable to premature solutioning. People arrive with preferred answers. A good agenda and interview guide should:

- Encourage fact collection before design decisions.
- Make decision rights and governance explicit topics.

- Produce tangible outputs (design principles, option set, next-step owners).
- Create psychological safety while preserving rigor.

A Level 1 chatbot is effective here because agendas and interview guides are structural artifacts. They do not require the model to know the client's operating model; they require it to propose a disciplined facilitation structure that humans can tailor.

**Governance reminders in this case.** The most common Level 1 risk is that the model will assume a "standard target operating model" and embed it implicitly. The control is to keep the agenda framed around elicitation and design criteria, not predetermined solutions. For example, the agenda should include segments like "agree design principles" and "define evaluation criteria," rather than "choose shared services" unless that direction is already decided and provided as a fact.

The interview guide should explicitly separate:

- **Observed facts:** what happens today, what decisions are made, what systems are used.
- **Pain points:** where delays occur, where accountability is unclear.
- **Desired outcomes:** what "better" would look like.
- **Constraints:** regulatory, system, talent, and change capacity constraints.

This structure helps prevent interviews from becoming purely opinion-driven and helps the team gather the inputs needed for later Level 2 reasoning.

**How to use these cases in practice.** These mini-cases are designed to be reused as training exercises. Teams can run them in a controlled environment, produce Draft Bundles, and practice reviewing outputs for invented facts and hidden assumptions. The point is not to perfect the drafts. The point is to build the habit of disciplined drafting.

Across all four cases, the consistent pattern is that Level 1 delivers value by accelerating structure and communication while protecting professional responsibility through explicit uncertainty and traceability. Market entry becomes a memo that frames verification. Cost transformation becomes a workplan and email that avoid false promises. Capital allocation becomes a board structure that flags missing data. Operating model redesign becomes a facilitation plan that elicits facts rather than assuming solutions. In every case, the human remains the owner of truth, and the chatbot remains a drafting tool.

## 1.7   Risks and controls taxonomy (Level 1)

Level 1 risk is often underestimated because the tasks look mundane. Drafting an email, summarizing a meeting, or producing a memo outline does not feel like a high-risk activity. In consulting and corporate strategy, however, those drafts are frequently the first formal record of a decision path. They shape expectations, create commitments, and travel across stakeholder networks with surprising speed. When a Level 1 chatbot is introduced into that environment, it changes the risk profile in a specific way: it makes *persuasive language cheap*. Cheap language is not inherently bad, but it lowers the friction that previously forced humans to pause, reconsider, and verify.

The purpose of a Level 1 risk taxonomy is not to dramatize AI. It is to name the predictable

failure modes that arise when fluent text is produced without embedded accountability. Once named, these failure modes can be controlled with a minimum set of practical measures that are easy to enforce under time pressure. A governance-first approach begins here: by treating failure as expected, and by designing workflows so that expected failures do not become professional incidents.

### 1.7.1 Failure modes you should expect

The following failure modes should not surprise you. They are not edge cases. They are normal behaviors of chatbots operating under incomplete inputs. The governance posture at Level 1 is therefore: assume these will occur unless you actively prevent them.

1. Hallucinated facts (numbers, events, "industry norms").
2. False confidence (polished language masking uncertainty).
3. Assumption drift (the model introduces new premises).
4. Scope creep (answering a different question than asked).
5. Confidentiality leakage (sensitive text pasted in).

**1) Hallucinated facts.** Hallucination is an overloaded term, but the practical issue is simple: the model will sometimes generate factual statements that were not provided, and it will state them as if they were known. In consulting and strategy, the most damaging hallucinations are those that look "reasonable":

- A plausible savings percentage in a cost program.
- A plausible market growth rate in a market entry memo.
- A plausible timeline for implementation.
- A plausible "industry standard" operating model.

These are dangerous because they can pass casual review. They are also dangerous because once they appear in a draft, they often persist. People copy and paste. Teams anchor on the first version. A number that begins as a guess can become treated as a fact by the third iteration.

Hallucinated facts are especially likely when the prompt implicitly asks for them. If you ask, "Draft a compelling market entry recommendation," the model will try to be compelling. If you ask for a "board pre-read" without providing financial context, the model may fill in what a board pre-read "normally" contains. This is why Level 1 governance insists on a "no new facts" rule and a structured output that makes evidence boundaries visible.

**2) False confidence.** Even when the model does not invent facts, it can still create a failure mode: confidence by tone. A chatbot can write in crisp, executive language. That language can make a tentative point feel settled. A paragraph that is actually an assumption can read like a conclusion. A list of open questions can disappear into "next steps" phrasing that implies certainty.

False confidence is not about lying. It is about style masking uncertainty. In professional settings, style is a credibility signal. The model is good at style. Therefore, without controls, the model will systematically increase the perceived credibility of unverified content.

This failure mode often appears in subtle linguistic choices:

- Replacing "might" with "will."
- Replacing "we should explore" with "we recommend."
- Replacing "preliminary" with "key finding."
- Replacing "unknown" with generic filler language that sounds resolved.

To control false confidence, the governance wrapper forces uncertainty into explicit fields: assumptions, open questions, and verification status. It also requires human review focused on tone and commitments, not only on grammar.

**3) Assumption drift.** Assumption drift is the model's tendency to introduce new premises that were not stated, often to make the draft coherent. Humans do this too, but humans are usually aware of the assumption they are making. The model is not. It will infer missing context and then treat that inferred context as part of the narrative.

In consulting, assumption drift is dangerous because it can change the direction of work. A memo draft that assumes "modest investment" means "distributor-only entry" can bias the team away from direct presence options. A workplan that assumes "no headcount reductions" means "only procurement savings" can narrow the hypothesis space prematurely. A stakeholder email that assumes "alignment already exists" can trigger political conflict when it does not.

Assumption drift is particularly likely when inputs are partial or ambiguous. It is also likely when the prompt asks for a complete artifact. Completeness is a trap: the model will try to complete the story. Level 1 governance counters this by demanding explicit assumptions and open questions. When assumptions must be listed, drift becomes visible. Reviewers can then decide whether to accept, revise, or remove them.

**4) Scope creep.** Scope creep is when the model answers a different question than the one asked. This can happen because the prompt was ambiguous, because the model defaults to common frameworks, or because the model "helpfully" expands the task. In consulting settings, scope creep can be operationally costly. If you ask for a meeting follow-up email and receive a full transformation roadmap, you may be tempted to use parts of it. If you ask for a one-page decision memo and receive a detailed recommendation with implied benchmarks, you may accidentally overreach in front of stakeholders.

Scope creep also creates governance risk. A prompt intended for drafting can drift into analysis-like output, which may violate Level 1 scope boundaries. For example, "create a risk section" can become "interpret regulatory requirements." "Draft a memo" can become "recommend a strategy." At Level 1, the control is not to prevent the model from producing extra content; it is to prevent the team from treating that extra content as legitimate. The output schema helps here: anything beyond provided facts must be captured as assumptions or open questions, not as conclusions.

**5) Confidentiality leakage.** The most serious Level 1 incidents are often not about wrong content; they are about wrong handling of content. Confidentiality leakage occurs when sensitive information is pasted into an unapproved interface, stored in an uncontrolled way, or shared beyond intended recipients. In consulting, this can include client names, deal terms, pricing, and internal performance data. In corporate strategy, it can include acquisition plans, restructuring details, and

competitive vulnerabilities.

This failure mode is primarily human, not model-driven. Under time pressure, people copy and paste. They paste too much. They paste without redacting. They paste into the easiest tool, not the approved tool. Governance controls must therefore be behavioral and enforceable: redaction defaults, minimum-necessary input, and approved environment rules.

Confidentiality leakage also includes a quieter form: outputs that inadvertently include sensitive information because the input contained it, and the output is then forwarded. Even if the original input was permissible, the output distribution might not be. This is why review gates and labeling matter: the person who sends the message must confirm it is appropriate for the audience.

### 1.7.2 Minimum control set (practical and enforceable)

A "risk taxonomy" is only valuable if it produces a control set that can actually be followed. At Level 1, the minimum controls must be:

- **Simple:** short enough to memorize.
- **Repeatable:** applicable across common drafting tasks.
- **Auditable:** capable of producing evidence of use.
- **Enforceable:** clear enough that managers can detect non-compliance.

The checklist below is the minimum standard for safe Level 1 operation. It is intentionally conservative. If an organization cannot follow these controls, it should restrict Level 1 usage to low-risk internal drafting until discipline is established.

---

**Checklist**

**Minimum controls (use in every engagement).**

1. **Schema requirement:** facts_provided / assumptions / open_questions / draft_output / verification_status.
2. **"No new facts" instruction:** the model must not add facts beyond provided inputs.
3. **Human review gate:** nothing leaves the team without reviewer initials and date.
4. **Versioning:** keep v0 model output and track edits (diff or change log).
5. **Confidentiality guardrails:** redaction + minimum-necessary input.

---

Each control maps directly to the failure modes.

**Control 1: Schema requirement.** This is the primary control against hallucinated facts, false confidence, and assumption drift. By forcing the model to list `facts_provided`, the team creates a visible evidence boundary. By forcing `assumptions`, the team makes hidden premises explicit. By forcing `open_questions`, the team prevents gap-filling from becoming assertion. By forcing `verification_status`, the team maintains the correct epistemic label: this is a draft.

Operationally, teams should treat schema compliance as a binary check. If the output does not include the schema, it is not an acceptable Level 1 output. This might sound strict, but strictness is what makes the control enforceable. If compliance is optional, it becomes rare.

**Control 2: "No new facts" instruction.** This is the direct antidote to hallucinated facts. The key is that the instruction must be explicit and repeated. In practice, it should appear in every prompt template, not only in policy documents. The model should be told:

- Do not invent numbers or benchmarks.
- Do not cite sources unless provided.
- If information is missing, list it as an open question.

The organization should also train reviewers to look for telltale signs of invented facts: phrases like "typically," "on average," "industry standard," and numeric ranges that were not provided. When detected, these must be removed or converted into open questions.

**Control 3: Human review gate.** This control addresses false confidence, scope creep, and confidentiality risk. Review is where professional judgment re-enters the loop. At Level 1, review should focus on a small set of questions:

a. Are all factual claims supported by provided facts?
b. Are assumptions explicitly listed and acceptable?
c. Are open questions clearly stated and routed?
d. Does the tone imply commitments we did not make?
e. Is the content appropriate for the audience (confidentiality)?

The requirement for initials and date is not bureaucratic. It is an accountability device. It ensures someone can be identified as the responsible reviewer. In fast-moving work, this can be as simple as a comment in the document: "Reviewed by X on date." What matters is that review is not implicit.

**Control 4: Versioning.** Versioning addresses assumption drift and provides governance trace. Keeping v0 (model output) and v1 (human edits) makes it possible to see what changed. It also supports learning. Over time, teams can identify recurring problems in v0 drafts and adjust prompt templates accordingly.

Versioning does not need to be complex. For emails, it can be as simple as saving the model output in the Draft Bundle folder before editing. For memos, it can be tracked in the document's revision history. The key is to preserve the provenance. If v0 disappears, you lose the ability to audit and to improve.

**Control 5: Confidentiality guardrails.** This is the control against confidentiality leakage. It has two parts: redaction and minimum-necessary input. Redaction reduces sensitivity. Minimum-necessary input reduces exposure. Together, they create a behavioral norm: do not paste more than you need.

Confidentiality guardrails also require clarity about approved environments. If the organization has specific approved tools or accounts, that must be communicated. If not, the governance-first posture is to assume that external-facing drafting should be performed with redacted inputs only.

**How to enforce the minimum controls.** Enforcement at Level 1 should be lightweight and cultural, not punitive. Practical enforcement mechanisms include:

- **Templates:** provide standard prompt templates that already include schema and "no new facts"

instructions.

- **Draft Bundle habit:** require saving the bundle for any deliverable that leaves the team or supports a decision.
- **Sampling:** periodically sample drafts for compliance and provide coaching.
- **Clear escalation:** define what to do when a draft touches regulated topics or high confidentiality.

The goal is to make safe behavior the default. Teams should not need to remember policy in the moment; the templates and workflow should make compliance easy.

**A final note on proportionality.** The minimum controls are a floor, not a ceiling. For some contexts, additional controls may be appropriate: second reviewer for board materials, restricted distribution lists, or formal sign-off steps. But at Level 1, the priority is to establish the baseline discipline that prevents predictable failures. Once that discipline is stable, the organization can decide whether and how to move to Level 2 reasoning patterns and beyond.

In governance-first adoption, Level 1 success is not measured by how many drafts the team produces. It is measured by whether the team can demonstrate a repeatable, defensible process: structured outputs, explicit uncertainty, human review, version history, and confidentiality discipline. If those are in place, Level 1 becomes a durable capability rather than a risky convenience.

## 1.8   Prompt patterns and exercises (copy/paste)

Level 1 succeeds or fails on repeatability. Most teams do not struggle because they lack "creativity" in prompting. They struggle because they do not have consistent, safe defaults under time pressure. When deadlines compress, people revert to convenience. They paste too much context. They ask for "a great memo." They accept outputs because they look polished. The prompt patterns in this section are designed to eliminate that drift by providing copy/paste templates that encode governance controls directly into the instruction text.

This section is deliberately practical. It treats prompts as operating procedures, not as clever tricks. A prompt is a control surface: it determines whether the model invents facts, whether uncertainty is surfaced, whether outputs are structured, and whether the artifact is reviewable. In a governance-first maturity ladder, prompt templates function like checklists in aviation: they standardize behavior so that predictable failures do not recur.

The templates below have three design principles:

1. **Strict boundaries:** "No new facts" is stated explicitly and repeatedly.
2. **Structured outputs:** facts, assumptions, and open questions are separated so reviewers can check quickly.
3. **Audit posture:** every output is labeled `"Not verified"` and produces a natural Draft Bundle.

You should treat these templates as the default starting point for Level 1. In most cases, you will modify only the `CONTEXT` and the `DELIVERABLE TYPE`. Resist the temptation to remove constraints for convenience. If you want more autonomy or more complex workflows, that is a sign you are moving into Level 2 or Level 3, which requires additional controls.

### 1.8.1 Prompt Template 1: Governance-first drafting wrapper (universal)

Template 1 is the universal wrapper for Level 1 drafting. It is appropriate for most common deliverables: stakeholder emails, meeting summaries, decision memos, workshop agendas, interview guides, and slide storylines (titles only). The template explicitly prohibits invented facts, enforces separation of facts and assumptions, and makes uncertainty visible.

Two practical notes before you use it:

- **Redact before pasting.** Replace client names and unique identifiers with placeholders. Keep the mapping offline.
- **Keep inputs sparse and factual.** The more you paste, the greater the confidentiality risk and the more likely the model is to embed irrelevant content.

---

**Prompt (Copy/Paste)**

```
ROLE: You are a drafting assistant for management consulting deliverables.
You must not invent facts, data, benchmarks, citations, or client details.
You must be explicit about uncertainty. You must separate facts from assumptions.


CONTEXT (provided by user):
[PASTE REDACTED NOTES HERE]


TASK:
Draft the deliverable requested below using ONLY the facts provided.
If information is missing, add it to open_questions (do not guess).
Use a professional consulting tone. Keep it concise and structured.


DELIVERABLE TYPE (choose one):
- stakeholder email
- one-page decision memo
- meeting summary / decision record
- slide storyline (titles only)
- workshop agenda
- interview guide


OUTPUT FORMAT (STRICT):
facts_provided: [bullet list]
assumptions: [bullet list]
open_questions: [bullet list]
draft_output: [the deliverable text]
verification_status: "Not verified"
questions_to_verify: [bullet list]
```

---

**How to use Template 1 well (practical guidance).** The template is strict by design. To get high-quality outputs within the constraints, you must provide the right kind of context. The

best context is:

- **Specific but bounded.** "We met with Head of Ops and CFO to discuss a four-week cost diagnostic" is good. "We are doing a transformation" is too vague.
- **Concrete facts over narratives.** Bullets are better than paragraphs. Facts are better than opinions.
- **Explicit uncertainty.** If you are unsure about a date or decision, mark it as "to confirm." This reduces the risk that the model upgrades it into a commitment.

**Review checklist for Template 1 outputs.** When the model returns the structured output, review it in this order:

1. **facts_provided:** Did the model restate your facts correctly? If not, fix inputs and re-run.
2. **assumptions:** Did the model add assumptions you did not intend? Remove or revise.
3. **open_questions:** Are the missing items surfaced? If the list is empty, be suspicious.
4. **draft_output:** Scan for invented numbers, "industry norms," or implied commitments.
5. **questions_to_verify:** Ensure the list is actionable and assign owners if possible.

**Common improvements (without breaking governance).** Teams often want the model to be more specific. The safe way to do that at Level 1 is not to let the model guess; it is to provide more verified facts. If you need the memo to mention a timeline, provide the timeline. If you need it to mention a KPI, provide the KPI. The model can be specific if you are specific. This is the right discipline: specificity follows verification.

### 1.8.2 Prompt Template 2: Tone variants without content drift

Template 2 is used when you already have a draft output (typically produced by Template 1 or by a human) and you want multiple tone variants for different audiences. This is an excellent Level 1 use case because it creates value without increasing epistemic risk *if* you prevent content drift.

The critical control is to treat meaning as invariant. Tone may change, structure may change, but the facts and caveats must not change. In practice, you should run Template 2 only after you have reviewed the base draft and removed any unverified claims. Otherwise, the model will simply produce three polished versions of the same mistake.

---

**Prompt (Copy/Paste)**

```
Using the draft_output below, rewrite it into 3 variants:
(1) neutral executive, (2) direct and concise, (3) diplomatic/stakeholder-friendly.
Do not add new facts. Keep meaning identical.


INPUT DRAFT:
[PASTE draft_output]


OUTPUT FORMAT:
variant_1:
variant_2:
variant_3:
verification_status: "Not verified"
```

---

**How to detect content drift.** Even with a "do not add new facts" instruction, models may change meaning subtly. Common drift patterns include:

- **Strengthening verbs:** "explore" becomes "pursue," "consider" becomes "recommend."
- **Dropping caveats:** "subject to validation" disappears in a more concise version.
- **Adding implied timing:** "soon" or "in the coming weeks" appears without being provided.
- **Introducing new rationale:** a "why" is added to make the message sound stronger.

To control this, compare variants against the base draft with a strict rule: if a variant includes any new claim, remove it or re-run with a stronger instruction such as "Do not add any new rationale, adjectives, or implied commitments; preserve all caveats verbatim."

**When to use which variant.**

- **Neutral executive:** updates to senior leadership; steering committee communications; board-adjacent notes.
- **Direct and concise:** internal coordination; action-focused messages; time-pressured stakeholders.
- **Diplomatic/stakeholder-friendly:** sensitive transformations; messages that require buy-in; communications to functions that may feel threatened.

This is not merely style. It is governance. The wrong tone can create political risk even when facts are correct. Template 2 helps teams tailor tone without creating new factual risk.

### 1.8.3   Exercise set (team training)

Level 1 adoption becomes durable when teams practice it as a craft. The exercises below are designed to build consistent habits: redaction, schema discipline, review gates, and Draft Bundle creation. These are not "AI prompts" exercises. They are *professional discipline* exercises conducted with AI assistance.

Each exercise should be run in small teams of two or three people. Require the team to produce a Draft Bundle for each exercise:

- Input snapshot (redacted)
- Prompt text used
- Output JSON (v0)
- Reviewer notes (what changed, what was verified, what remains open)
- Version history (v0, v1, v2)

The goal is to train both the person who drafts and the person who reviews. In real work, failures often occur because the reviewer assumes someone else checked. These exercises build a shared standard.

### Exercise 1: Notes-to-email with commitment control (Market entry follow-up).

1. Provide a short set of redacted meeting notes for a market entry discussion:

   - CEO wants a decision framework in two weeks.
   - Country B is attractive due to mining sector growth (claimed by stakeholder, unverified).
   - Two options: distributor vs direct presence.
   - Compliance and service capability are key risks.
   - Next meeting scheduled "next week" (date not confirmed).

2. Use Template 1 to draft a stakeholder follow-up email.
3. Reviewer task: remove any implied commitments and ensure all timing language is either confirmed or labeled as "to confirm."
4. Deliverable: v2 email draft + a list of open questions routed to owners.

**Learning objective:** Train the team to detect implicit commitments introduced by the model and to enforce cautious language.

### Exercise 2: Decision record with strict decision rule (Steering meeting).

1. Provide a pseudo-transcript (bullets by speaker role) including tentative statements and one explicit decision.
2. Use Template 1 with "meeting summary / decision record."
3. Reviewer task: ensure only explicitly confirmed items are labeled as decisions; everything else becomes discussion points or open items.
4. Deliverable: a final decision record and an action tracker table (Owner: TBD allowed).

**Learning objective:** Prevent the model from upgrading discussion into decision; practice writing "Owner: TBD" instead of guessing.

### Exercise 3: Outline-to-memo with assumption surfacing (Cost diagnostic memo).

1. Provide an outline for a four-week cost diagnostic, including constraints: "no headcount decisions in first four weeks" and "savings ranges only after validation."
2. Use Template 1 with "one-page decision memo."
3. Reviewer task: ensure the memo includes an explicit assumptions section and a clear verification plan; remove any invented savings benchmarks.
4. Deliverable: a memo draft that a CFO could read without being misled.

**Learning objective:** Train teams to insist that uncertainty and validation gates are visible in executive writing.

**Exercise 4: Board pre-read outline with placeholders (Capital allocation).**

1. Provide a fact pattern with three uses of capital and missing forecast data.
2. Use Template 1 with "slide storyline (titles only)" or "one-page decision memo," depending on preference.
3. Reviewer task: enforce placeholders for missing data ("Insert verified forecast," "Insert covenant headroom"); remove any implied recommendation.
4. Deliverable: an outline or storyline that is decision-driven and explicitly data-dependent.

**Learning objective:** Practice producing high-stakes structures without overclaiming and without invented numbers.

**Exercise 5: Workshop agenda + interview guide (Operating model redesign).**

1. Provide workshop constraints: 3 hours, change fatigue, cross-functional participants.
2. Use Template 1 to produce a workshop agenda and an interview guide.
3. Reviewer task: ensure agenda segments are designed for elicitation and alignment, not predetermined solutions; ensure interview questions separate facts from opinions.
4. Deliverable: agenda + guide + a short list of "governance reminders" to read at the start of the workshop.

**Learning objective:** Train teams to use Level 1 for facilitation structure while preserving neutrality and avoiding solution bias.

**Exercise 6: Tone variants with invariants (Stakeholder sensitivity).**

1. Take the v2 email from Exercise 1 or 3 as the base draft.
2. Use Template 2 to produce three tone variants.
3. Reviewer task: compare each variant to the base draft and highlight any content drift; correct and re-run if needed.
4. Deliverable: three variants plus a one-paragraph note explaining which audience each is for.

**Learning objective:** Build the discipline of treating tone rewriting as a controlled transformation, not a rewriting free-for-all.

**How to run these exercises as a team standard.** To institutionalize Level 1 discipline, run these exercises as part of onboarding for new team members or as part of periodic capability refresh. The goal is to create a shared understanding that:

- AI drafts are v0, not final.
- Facts must be traceable; assumptions must be owned.
- Uncertainty must be explicit, not hidden.
- Review and versioning are mandatory, not optional.

If the team can execute these exercises consistently, it is ready to use Level 1 in real engagements. If the team cannot, the organization should restrict AI usage to very low-risk internal drafting until

the governance habits are stable.

These prompt patterns and exercises complete the Level 1 objective: make chatbot drafting valuable, repeatable, and defensible. They also set up the natural transition to Level 2. At Level 2, the same discipline will be applied not just to drafting, but to explicit reasoning: issue trees, tradeoffs, and assumption registers that block downstream steps. But the core posture remains the same: humans own truth, models support structure, and every output is governed until verified.

## 1.9 Conclusion and transition to Level 2 (Reasoners)

Level 1 is the moment when an organization learns a new professional discipline: how to work in a world where polished language is inexpensive. That may sound like a minor shift, but it is not. In consulting and corporate strategy, credibility is built through the controlled relationship between claims and evidence. Historically, the cost of producing polished text created friction that indirectly protected that relationship. When a chatbot can draft a crisp memo in seconds, that friction disappears. What replaces it must be explicit: a governance wrapper that preserves epistemic discipline and professional accountability.

This chapter has treated Level 1 as exactly what it is: a drafting accelerator with real operational value and predictable risks. It has also treated Level 1 as a foundation. The habits formed here determine whether the organization's later AI maturity becomes defensible capability or scalable failure. Teams that treat Level 1 as "just a writing tool" often discover, too late, that writing is where assumptions propagate, commitments are created, and reputational risk is born. Teams that treat Level 1 as governed drafting discover something better: that structure and transparency improve communication quality, reduce misalignment, and strengthen trust.

The transition to Level 2 is not a jump in "intelligence." It is a shift in how the tool is used. Level 1 focuses on expressing what the team already knows. Level 2 focuses on structuring what the team does not yet know: the reasoning pathway from facts to decisions, the alternatives available, the tradeoffs implied, and the verification plan required. That shift can create substantial value, but it also increases risk because the model begins to influence framing, not merely phrasing. Therefore governance must strengthen accordingly.

### 1.9.1 Summary of main takeaways

The most important outcome of Level 1 is not a library of prompts. It is a shared understanding of what a chatbot is *in professional work* and how it must be governed.

**First, Level 1 is drafting acceleration, not analysis and not verification.** A chatbot can transform rough notes into a structured email, memo, decision record, or storyline. It can enforce consistency in sections and tone. It can create checklists and agendas from your inputs. But it cannot be treated as a source of truth. It does not validate facts. It does not "know" your client context. It does not bear accountability for errors. When the model produces a credible-sounding statement, that credibility comes from style, not from evidence. This is why Level 1 must be framed

explicitly as drafting support.

**Second, the biggest Level 1 risk is persuasive error, not visible nonsense.** In practice, the most dangerous outputs are not obviously wrong. They are subtly wrong, plausibly wrong, and therefore easy to ship. A model can invent a benchmark that looks reasonable, strengthen tentative language into commitment language, or smooth away caveats that were the only honest part of the draft. In consulting and strategy, these errors can create expectations that are costly to unwind. The governance-first posture therefore treats polished language as a risk signal: the cleaner the draft, the more carefully it must be reviewed for factual grounding and implied certainty.

**Third, the foundational control is epistemic separation: facts are not assumptions.** This chapter has repeatedly insisted on a schema that separates `facts_provided`, `assumptions`, and `open_questions`. This separation is not cosmetic. It is the minimum mechanism that prevents hallucinated facts, assumption drift, and scope creep from quietly entering deliverables. It also makes review efficient. Reviewers can scan facts, challenge assumptions, and route open questions without reading every sentence as a potential risk. In other words, the schema turns drafting into a controlled process.

**Fourth, Level 1 governance is light but non-negotiable.** The minimum wrapper is simple:

- A standard input/output schema.
- A "no new facts" instruction.
- Redaction and minimum-necessary input.
- A human review gate and explicit sign-off.
- Version history and Draft Bundle preservation.

None of these controls are exotic. They are consistent with how professional work should already be governed. The difference is that AI makes it easy to skip them. Therefore, Level 1 success requires making them explicit and habitual.

**Fifth, Level 1 value scales through repeatable workflow patterns.** Notes-to-email, transcript-to-decision record, outline-to-memo, and storyline drafting are recurring tasks in consulting and corporate strategy. When teams adopt standardized Level 1 patterns, they reduce variability, improve clarity, and build internal IP. The value is not merely speed; it is consistency. Consistency is what allows teams to coordinate better and to defend process when challenged.

**Sixth, Level 1 is the training ground for organizational maturity.** By requiring Draft Bundles, prompt logs, reviewer notes, and version history, Level 1 introduces the core disciplines that later enable agentic workflows and enterprise governance. If those disciplines cannot be executed at Level 1, it is not safe to scale to higher levels. This is a critical point: the maturity ladder is not a ladder of ambition; it is a ladder of governability.

Taken together, these takeaways define what Level 1 is meant to accomplish: a safe, repeatable drafting capability that makes consulting and strategy work faster without making it less defensible.

### 1.9.2  What comes next (preview of Level 2)

Level 2 introduces a new use of generative AI: structured reasoning support. The shift is subtle but profound. At Level 1, the model helps you express and package what you already have. At Level 2, the model helps you structure the path from what you have to what you need: the decomposition of a problem, the articulation of alternatives, the mapping of tradeoffs, and the identification of the weakest assumptions that could break a recommendation.

This is where many teams are tempted to anthropomorphize the tool. A model that can generate an issue tree can feel like an analyst. A model that can list alternatives can feel like strategic insight. A model that can propose tradeoffs can feel like judgment. The governance-first posture insists on a different interpretation: Level 2 is a reasoning scaffold. It is a tool for making human reasoning more explicit, more checkable, and more complete. It does not replace judgment; it makes judgment auditable.

The core capabilities introduced in Level 2 typically include:

**Issue trees and structured decomposition.** Instead of starting with a narrative draft, Level 2 starts with a map of the problem. For a market entry decision, the model can help lay out the dimensions that must be evaluated: market attractiveness, competitive dynamics, route-to-market options, capability gaps, regulatory constraints, and financial implications. The value is not that the model "knows" these categories. The value is that it forces the team to see the full space and to avoid missing a critical branch.

**Alternatives and option sets.** Level 2 formalizes the discipline of alternatives. In consulting and strategy, poor decisions often arise from false binaries: only one option feels viable. A reasoning scaffold can help teams generate a structured option set, including conservative and aggressive variants. Governance requires that these alternatives be labeled as exploratory, not as recommendations, and that the assumptions behind each alternative be explicit.

**Tradeoffs and decision criteria.** Level 2 introduces explicit tradeoff framing: what do we gain and what do we give up under each option? What criteria does leadership care about? Speed, risk, control, investment, strategic flexibility, reputational impact? By forcing criteria to be written down, Level 2 reduces the risk of hidden value judgments.

**Assumption registers and weakest-link checks.** This is the major governance enhancement at Level 2. Instead of treating assumptions as a list attached to a draft, Level 2 treats assumptions as managed objects. They are tracked, prioritized, and used to gate downstream outputs. If a key assumption is unverified, the workflow must not allow a draft to present conclusions that depend on it. This is how governance strengthens: the organization moves from "we listed assumptions" to "assumptions control what we are allowed to claim."

**Verification planning.** Level 2 makes verification a first-class deliverable. Instead of merely listing open questions, the team produces a verification plan: what data is needed, where it will come from, who owns it, what is the timeline, and what thresholds would change the decision. This is the bridge between reasoning and execution. It also makes the work auditable: a reviewer can see whether the team verified what it said it would verify.

Because Level 2 influences framing and reasoning, it introduces new risks:

- The model may bias the framing toward conventional frameworks.
- The model may omit critical branches, giving false completeness.
- The model may propose plausible tradeoffs that do not fit the specific context.
- The model may implicitly push toward a recommendation without adequate verification.

Therefore governance must strengthen. Level 2 will require clearer separation between facts, assumptions, and reasoning steps; stricter assumption gating; more explicit reviewer roles; and more robust artifact logging for reasoning outputs. In short, Level 2 is not "Level 1 plus smarter prompts." It is a different operating mode with higher stakes and therefore higher control needs.

The transition is best understood as a shift from *drafting artifacts* to *drafting reasoning structures.* Level 1 produces emails and memos. Level 2 produces issue maps, decision matrices, assumption registers, and verification plans that then feed into those emails and memos. If Level 1 improves communication, Level 2 improves the transparency of thinking. But transparency of thinking is only valuable when it is coupled with governance that prevents speculative reasoning from being mistaken for validated analysis.

---

**Artifact (Save This)**

**Level 1 exit criteria (ready to move to Level 2).**

1. The team consistently uses facts/assumptions/open_questions in outputs.
2. Draft bundles are saved with version history and reviewer notes.
3. No deliverables are sent externally without verification and sign-off.
4. The team can demonstrate a repeatable prompt pattern library (internal IP).

---

These exit criteria are not bureaucratic hurdles. They are safety conditions. They ensure that the organization has proven it can use AI for drafting without losing control of truth, accountability, and confidentiality. Once these criteria are met, Level 2 becomes a responsible next step: adding structured reasoning in a way that remains governed, auditable, and owned by humans.

In the next chapter, Level 2 will formalize reasoning patterns that support consulting and corporate strategy work: issue trees, alternative generation, tradeoff mapping, assumption registers that gate outputs, and verification plans that turn uncertainty into an executable workplan. The tool will become more helpful, but the obligations will become stronger. That is the discipline of the maturity ladder: capability increases, governance increases, and professional judgment remains the human responsibility at every step.

# Bibliography

[1] National Institute of Standards and Technology (NIST). *Artificial Intelligence Risk Management Framework (AI RMF 1.0).* NIST AI 100-1, January 2023.

[2] ISO/IEC. *ISO/IEC 42001:2023 Information technology — Artificial intelligence — Management system.* International Organization for Standardization, 2023.

[3] ISO/IEC. *ISO/IEC 23894:2023 Information technology — Artificial intelligence — Guidance on risk management.* International Organization for Standardization, 2023.

[4] ISO. *ISO 31000:2018 Risk management — Guidelines.* International Organization for Standardization, 2018.

[5] Organisation for Economic Co-operation and Development (OECD). *Recommendation of the Council on Artificial Intelligence.* OECD Legal Instrument OECD/LEGAL/0449, adopted May 2019.

[6] European Union. *Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act).* Official Journal of the European Union, 2024.

[7] UK Information Commissioner's Office (ICO). *Guidance on AI and data protection.* Updated March 2023.

[8] Association for Computing Machinery (ACM). *ACM Code of Ethics and Professional Conduct.* Adopted June 22, 2018.

# Chapter 2

# Reasoners

**Abstract.** Level 2 (Reasoners) introduces disciplined structure to the use of Generative AI in management consulting and corporate strategy. Where Level 1 accelerates drafting, Level 2 externalizes thinking into inspectable artifacts: issue maps, alternatives, tradeoff tables, assumption registers, and verification plans. The objective is not to outsource judgment, but to reduce "persuasive error" by forcing explicit separation between what is known, what is assumed, and what must be verified before any decision or client-facing recommendation is made. In this chapter, the model is treated as a reasoning scaffold that helps teams decompose ambiguous questions, enumerate options (including the status quo), and surface second-order effects, constraints, and dependencies. The governance posture tightens accordingly: the model must not invent facts, benchmarks, or citations; it must not recommend or select an option; and every assumption introduced must be logged with an owner role and a concrete test. Outputs remain drafts with verification_status set to `"Not verified"` until reviewed and signed off by a responsible professional. The chapter provides reusable workflow patterns (A–D) that move from question framing to weakest-link assumption testing, plus recurring mini-cases (market entry, cost transformation, capital allocation, and operating model redesign) that demonstrate how structured reasoning artifacts support executive discussion, not decision automation. Readers finish with a repeatable "Reasoning Bundle" standard that strengthens defensibility, traceability, and review discipline as AI capability increases.

---

**Artifact (Save This)**

**Scope disclaimer (required).** This chapter is an educational governance-first guide for management consulting and corporate strategy work. Outputs produced with AI are drafts only and must be reviewed by a qualified professional. The model may be persuasive and wrong. Do not treat reasoning artifacts as correct, complete, or client-ready without verification and sign-off.

**Level alignment.** Level 2 uses a single model to support *structured reasoning*: issue maps, alternatives, tradeoffs, assumption registers, and verification planning. No autonomous agents, no tool use, no independent research is assumed. Human judgment remains fully responsible.

---

## 2.1 Chapter overview: Level 2 in the maturity ladder

Level 2 exists because consulting work does not fail only when a document is poorly written; it fails when reasoning is implicit, when assumptions are smuggled in as if they were facts, and when a polished narrative creates the illusion that the team has done the hard work of thinking. Level 1 (Chatbots) is valuable precisely because it improves speed and clarity: it turns rough notes into a usable email, converts a messy meeting into a decision record, and helps impose structure on a memo that otherwise would remain fragmented. But Level 1 also amplifies a well-known professional hazard: a coherent narrative can be persuasive even when it is wrong. In consulting and corporate strategy, this is not a cosmetic problem. The artifact is often the mechanism by which a decision travels through an organization. If the artifact is persuasive, it can move faster than the underlying verification. Level 2 is the deliberate answer to that risk. It makes thinking visible as structured

objects that can be challenged and verified: issue maps, alternatives, tradeoffs, assumption registers, and verification plans. In short, Level 2 moves the model from "drafting assistant" to "reasoning scaffold," and it raises the governance bar accordingly.

The maturity ladder is not a promise that higher levels are better in every context. It is a control framework. Each higher level introduces new capabilities that can create new value, but also introduces new failure modes that are harder to detect and more expensive to unwind. The ladder helps teams answer a practical question that is too often ignored: "Given the task we are performing, what is the minimum capability we need, and what is the minimum governance we must implement to use that capability safely?" Level 2 is the point in the ladder where the work product stops being "a draft" and starts becoming "a traceable reasoning bundle." That is also where the risk of decision laundering becomes real: outputs can look like analysis. This chapter therefore makes the governance posture explicit: at Level 2, the model is prohibited from recommending, choosing, or deciding. It is permitted to structure, compare, and stress-test, but not to conclude. It is permitted to identify what would need to be true for an option to work, but not to claim that it is true. It is permitted to propose verification questions, but not to answer them without evidence. Those boundaries are not rhetorical. They are the difference between a defensible workflow and a workflow that quietly replaces professional judgment with fluent text.

### 2.1.1 The five levels (why this ladder exists)

The five-level ladder exists to align capability with accountability. In most organizations, "using AI" is treated as a single category, as if drafting a meeting recap and running a multi-step workflow with automated retrieval and delegated subtasks were essentially the same thing. They are not. They differ in how much autonomy is implied, how much information the system touches, how often the system can be wrong without being noticed, and how quickly errors can propagate into client decisions or board-level commitments. A ladder forces a simple discipline: name the capability, name the risk profile, and name the minimum controls.

Level 1, Chatbots, is the entry point. The system is used as a drafting accelerator. It produces text artifacts: emails, memos, outlines, meeting notes, and first-pass checklists. The core value is speed and formatting; the core risk is persuasive error and the accidental introduction of facts the team did not supply. The governance at Level 1 is therefore centered on input discipline (minimum-necessary, redacted facts) and output discipline (facts vs assumptions vs open questions, explicit "Not verified," reviewer sign-off). In many firms, Level 1 is already "enough" for most day-to-day work, because much of consulting is communication. But Level 1 becomes dangerous when teams implicitly treat the output as analysis rather than as a draft.

Level 2, Reasoners, adds explicit reasoning structure. The model is used to make thinking inspectable rather than to make prose prettier. It decomposes questions into issue maps, enumerates alternatives, articulates tradeoffs, and stresses assumptions. The core value is cognitive discipline: it makes it harder to hide uncertainty behind narrative flow. The core risk is false rigor: the structure itself can create unjustified confidence. This is the paradox of Level 2. A good issue tree can look

like an answer. A tradeoff table can feel like due diligence. An assumption register can appear to be verification. None of those things are verification. They are scaffolds for verification. Therefore the governance at Level 2 tightens around three non-negotiables: (1) strict separation of facts and assumptions, (2) a no-recommendation rule, and (3) a verification plan that assigns ownership for checking what matters. This chapter is about building that discipline.

Level 3, Agents, introduces multi-step workflows. At this level the work is not just "one prompt produces one artifact." The system may execute a sequence: produce an issue map, generate interview questions, collect inputs from multiple stakeholders, synthesize results, run a QA pass, and package deliverables with logs. Even if you use only one model, the workflow itself behaves like an agentic process because it performs multiple dependent steps. The value is throughput and repeatability. The risk is compound error and process opacity: if steps are not logged and gated, small mistakes can cascade into large outputs. Governance therefore shifts from "output disclaimer" to "process governance": checkpoints, separation of duties, immutable logs, and kill-switch criteria. Level 3 is where auditability becomes a design constraint, not a compliance afterthought.

Level 4, Innovators, is about reusable assets. Here the firm begins to treat prompts, templates, evaluation rubrics, and workflow orchestrators as intellectual property that can be versioned, tested, and deployed across teams. The value is consistency and leverage. The risk is systemic: a flawed asset can scale quickly across projects, and governance failures become portfolio-level failures. Controls therefore resemble software governance: controlled releases, regression testing, documentation of scope and limitations, and a change log. It is also where training, onboarding, and supervision frameworks become part of the asset, because the way the asset is used determines its risk profile.

Level 5, Organizations, is firm-level governance. The firm operates a system of record: intake classification, data controls, model and prompt inventories, evaluation harnesses, audit trails, and escalation paths. Human accountability is explicit and enforced. At this level, the organization acknowledges that "AI usage" is not a set of individual choices; it is an operating model. The value is scale with defensibility. The risk is organizational: improper use can create legal, reputational, and client trust damage. Controls must therefore cover policy, process, training, tooling, and monitoring, with governance that is visible to leadership.

The ladder matters in consulting and corporate strategy because consulting outputs are often decision triggers. A one-page memo can launch a cost transformation. A slide storyline can steer capital allocation. A set of interview questions can shape what the team believes is true. If AI is used without a ladder, teams tend to slip into one of two bad patterns. The first is overconfidence: assuming that any AI use is "advanced," and therefore treating the outputs as special. The second is under-governance: assuming that AI use is "just drafting," and therefore failing to apply controls when the output is actually being used as analysis. The ladder creates the missing middle ground. It gives teams a vocabulary to say: "This task is Level 2. We are using reasoning scaffolds, not agents. Therefore we need assumption registers and verification planning, but we do not need a workflow orchestrator or automated retrieval." Or: "This task is Level 3. We are running a multi-step process; therefore we need checkpoints and immutable logs." The ladder prevents capability drift: the quiet

escalation from a chatbot draft to an agentic system without governance.

Level 2 is particularly important because it is the first level where outputs can plausibly be mistaken for "analysis." This is the level where discipline must be taught, practiced, and enforced. In real engagements, the first time a team produces a clean issue map and a table of alternatives, someone in the room will be tempted to treat it as an answer. Level 2 governance exists to prevent that temptation from becoming a habit.

### 2.1.2 Why Level 1 is not enough

Level 1 is not enough because consulting work is not primarily a writing contest. The deliverable is a vehicle for reasoning, and reasoning has failure modes that good writing can conceal. A chatbot can be a superb editor. It can take scattered thoughts and turn them into a narrative with clean headings, crisp bullets, and a confident tone. Those are useful outputs. But that same fluency makes it too easy to ship an artifact before the team has done the hard work of problem structuring and verification.

The first failure mode at Level 1 is persuasive prose. The model produces text that sounds like it belongs in an executive memo: balanced tone, well-shaped paragraphs, and decisive language. But persuasiveness is not evidence. When a team is under time pressure, a fluent draft can reduce the felt urgency to validate claims. Instead of acting as a mirror of what the team knows, the draft begins to feel like a completion of the work. The organization experiences "closure," even when the underlying question remains open. In corporate settings, closure is not neutral. It changes behavior. It shifts meeting agendas from exploration to execution. It can trigger communications to stakeholders that are difficult to reverse. Level 2 exists to keep the organization in the correct mode: exploration until verification, then decision.

The second failure mode is hidden assumptions. At Level 1, a chatbot is often asked to "fill in the gaps" to make the draft read smoothly. Even when you instruct the model not to invent facts, it can introduce soft assumptions that pass unnoticed: implied timelines, typical margin ranges, presumed stakeholder alignment, assumed feasibility of hiring or automation, or unspoken legal and regulatory constraints. These assumptions are not necessarily malicious; they are the model doing what language models do: completing patterns. But in consulting, assumptions are not harmless. They become the premises of downstream analysis. If assumptions are not surfaced and owned, they become invisible drivers of the recommendation that later emerges. A Level 1 draft can therefore become a Trojan horse: a clean narrative that quietly embeds untested premises.

The third failure mode is false completeness. Many consulting questions are underdetermined. "Should we enter this market?" "How much cost can we take out?" "Should we consolidate functions?" These questions cannot be answered without defining scope, objective function, constraints, risk appetite, time horizon, and decision authority. A Level 1 system will happily draft an answer-shaped memo even when those elements are missing. The memo may include "risks" and "next steps," which makes it appear responsibly cautious, but it can still be structurally wrong because the underlying decision frame is absent. Level 2 responds by forcing the decision frame to be explicit:

what decision is being made, by whom, by when, against what criteria.

The fourth failure mode is the collapse of uncertainty. In strategy work, uncertainty is not an inconvenience; it is often the central fact. Market data may be incomplete. Stakeholder incentives may conflict. Execution capacity may be unknown. A Level 1 draft tends to smooth uncertainty into generalities. It may replace "we do not know" with "it depends," which is not the same thing. It may replace "we should verify" with "we should consider," which weakens accountability. Over time, teams can become accustomed to a style of output that feels responsible but is not operationally enforceable. Level 2 counteracts this by turning uncertainty into explicit objects: open questions, verification tasks, and weakest-link assumptions.

A fifth failure mode is decision laundering. This is subtle and common. A leader has a preferred option. The team drafts a memo. The chatbot helps craft a narrative that supports the option, with plausible-sounding framing and a set of "tradeoffs" that make the option look balanced. The memo is then presented as analysis. The organization tells itself that it made a reasoned decision, but in reality it used the drafting tool to manufacture legitimacy. Level 1 cannot prevent this. Level 2 can mitigate it by forcing alternatives to be explicit, by requiring the status quo to be stated, and by demanding that the assumptions that make the preferred option viable be named and tested. Level 2 does not remove politics from decision-making. It does something more modest and more valuable: it makes the logic legible, so that politics cannot hide behind fluent prose.

Finally, Level 1 is not enough because it does not naturally produce artifacts that are durable under scrutiny. A consulting team may be asked weeks later: "Why did we think this?" "What did we assume?" "What did we verify?" "What evidence did we have at the time?" A Level 1 workflow often cannot answer those questions without reconstructing history from memory and scattered notes. This is not a theoretical problem. In real organizations, decisions are revisited, sometimes under stress. When results disappoint, leadership wants traceability. If you cannot show your reasoning and verification posture at the time, you lose credibility. Level 2 is the simplest step that makes reasoning reconstructable.

Level 2 therefore exists not because Level 1 is "bad," but because Level 1 is insufficient for the types of accountability that strategy work demands. Level 1 is drafting discipline. Level 2 is thinking discipline. And in consulting, thinking discipline is the difference between a defensible recommendation and an elegant guess.

### 2.1.3 What changes at Level 2 (and what does not)

The most important change at Level 2 is that the unit of work is no longer a document draft. The unit of work becomes a reasoning bundle: a set of structured artifacts that represent how the team is thinking. The memo or slide deck may still exist, but it is downstream of the reasoning objects. In Level 1, structure is applied to writing. In Level 2, structure is applied to the decision itself.

The first concrete change is issue mapping. Level 2 begins by decomposing an ambiguous question into sub-questions that can be investigated. This is not a mechanical exercise. It forces the team to decide what "the question" actually is. For example, "Should we enter Market X?"

might decompose into: attractiveness, right-to-win, economics, risks, and timing. Each branch then reveals missing information and implicit assumptions. The value of the issue map is not that it is perfect; it is that it makes the scope of reasoning explicit and therefore challengeable. The governance implication is that the issue map must be treated as a draft artifact. It can be wrong. It must be reviewed. It must not be used as proof of completeness.

The second change is alternatives discipline. Level 2 requires the explicit enumeration of options before the team converges. This is a control against premature convergence and confirmation bias. A model can help generate alternatives, but governance requires that alternatives be framed as options, not as recommendations. The model must not "pick" the best one. It must present a set, including the status quo or a "do nothing" option when appropriate. The human team owns the decision about which alternatives are real contenders and which are straw men.

The third change is tradeoff mapping. At Level 2, tradeoffs are not decorative. They are structured representations of how each alternative performs against criteria, constraints, and risks. Tradeoffs must include second-order effects and dependencies: what must be true for the alternative to succeed, what the organization would have to give up, and what risks are irreversible. This makes the reasoning more robust, but also more dangerous if misused. A tradeoff table can be mistaken for a quantified business case even when it is purely qualitative. Governance therefore requires explicit labeling: these tradeoffs are hypotheses, not verified conclusions, unless supported by evidence.

The fourth change is the assumption register. This is the signature artifact of Level 2. At Level 1, assumptions are often an afterthought, listed briefly to appear cautious. At Level 2, assumptions are treated as first-class objects with ownership and testability. An assumption register is not a list of vague statements like "assume stable demand." It is a structured log: each assumption, why it matters, how it could fail, who owns it, and what evidence would test it. The assumption register is the primary mechanism for making reasoning auditable. It is also the main defense against the model's tendency to introduce new premises. If an assumption is introduced, it must be logged. If it is logged, it can be challenged. If it can be challenged, it can be verified or rejected.

The fifth change is verification planning. Level 2 ends not with an answer, but with a plan to reduce uncertainty. The model can propose verification questions, but cannot answer them. The plan assigns tasks: what evidence is required, who will gather it, and what decision threshold or trigger would cause the team to change its mind. Verification planning is what prevents Level 2 from becoming a sophisticated form of speculation. Without verification planning, Level 2 is just structured storytelling. With verification planning, Level 2 becomes a disciplined pre-decision process.

These changes come with a deeper governance posture: Level 2 must prevent decision laundering. The purpose of a reasoning bundle is to make it harder to smuggle a preferred conclusion into a professional-looking artifact. Therefore, Level 2 includes explicit controls such as "weakest-link assumptions," which forces the team to identify the small set of assumptions that, if wrong, would flip the decision. This is a powerful discipline because it converts broad uncertainty into prioritized verification.

Just as important is what does not change. Level 2 does not verify facts. It does not independent research. It does not provide authoritative benchmarks, market sizing, or competitor intelligence unless those inputs are supplied and independently validated. Level 2 does not recommend. It does not choose. It does not decide. It can compare options and articulate how they differ, but it cannot tell the organization what to do. If a team wants a recommendation, that recommendation must be owned by a human professional, supported by evidence, and documented as judgment. The model can support the preparation of that judgment by making the logic explicit, but it cannot replace it.

Level 2 also does not remove the need for human context. Strategy decisions are embedded in organizations: politics, incentives, timing, and relationships matter. A model can help surface stakeholders and constraints as structured items, but it cannot know the actual dynamics unless you provide them, and even then it can only mirror the inputs. Governance therefore requires a human owner of the decision frame and a reviewer who challenges the reasoning artifacts for completeness, bias, and missing context.

Finally, Level 2 does not remove confidentiality risk. In fact, it can increase it, because teams are tempted to provide more context "to improve reasoning." This chapter therefore treats redaction and minimum-necessary input as ongoing requirements, not Level 1-only advice. The more strategic the question, the more sensitive the inputs often are: executive intent, negotiation posture, internal performance gaps, and restructuring plans. Level 2 must be practiced within approved environments and with disciplined redaction, because a better reasoning scaffold is not worth a breach of trust.

The result of these changes is a different kind of professional output. In Level 1, the model produces a draft artifact that a human edits. In Level 2, the model helps produce a reasoning package that a human uses to decide what to verify, what to debate, and what to escalate. The "deliverable" is not the issue map alone, nor the alternatives list, nor the tradeoff table. The deliverable is the ability to show, after the fact, that the team's thinking was structured, that assumptions were explicit, that uncertainty was acknowledged, and that verification was planned. That is the purpose of Level 2 in the maturity ladder: not more intelligence, but more discipline.

## 2.2 Mental model: what a Reasoner is in consulting work

Level 2 is where many teams first get into trouble precisely because the outputs become *convincing* in a way that Level 1 outputs usually are not. A chatbot draft can be dismissed as "just wording." A reasoning bundle, however, looks like analysis. It has a tree, a table, an assumption register, and a plan. It looks like work. It looks like rigor. And because it looks like rigor, it can short-circuit the two habits that make consulting trustworthy: professional skepticism and explicit verification. The mental model for Level 2 must therefore be conservative, operational, and governance-first. A Reasoner is not "a smarter chatbot." A Reasoner is a disciplined tool for transforming ambiguous executive questions into inspectable reasoning objects. It is a *scaffold* that supports human thinking and a *structured challenger* that makes it harder for humans to hide behind narrative flow.

The correct way to think about a Reasoner is to treat it like a junior analyst who is very fast at structuring, very fluent in frameworks, and very capable of producing plausible-sounding logic, but who has two critical limitations: it does not know what is true unless you provide and verify it, and it cannot be accountable for decisions. If you would not let a junior analyst sign a recommendation memo without a manager review and evidence, you should not let a Reasoner's output travel as if it were a conclusion. Level 2 is about making the junior analyst's work product better: clearer logic, explicit assumptions, and a verification plan. But it is still the manager's job to validate, choose, and own.

The chapter's mental model is built around a simple transformation: Level 1 converts *notes into narrative.* Level 2 converts *narrative into objects.* Those objects—issue maps, alternatives, tradeoffs, assumption registers, verification plans—are not "analysis" by themselves. They are intermediate artifacts that allow humans to do analysis more reliably. They are designed to be examined, criticized, and revised. The output is incomplete by design because completeness cannot be declared without verification. When a Reasoner is used correctly, it produces a visible structure that invites scrutiny. When it is used incorrectly, it produces a visible structure that suppresses scrutiny.

### 2.2.1 Useful abstraction

A useful abstraction for a Reasoner in consulting work is: *a reasoning scaffold plus a structured challenger.* The scaffold is the part that helps you build: it gives you a disciplined way to represent a problem, to define its boundaries, to enumerate options, and to organize uncertainty. The challenger is the part that helps you test: it identifies weak links, surfaces hidden assumptions, and asks the questions your team might be too busy—or too biased—to ask.

**The Reasoning Scaffold.** Consulting problems are usually under-specified. An executive question arrives as a sentence: "Should we enter this market?" "How much cost can we remove?" "What operating model should we adopt?" But the work is not in answering the sentence; it is in defining what the sentence *means.* A Reasoner assists by turning the sentence into a decision frame and then into inspectable components. That conversion is the core value. The model takes an ambiguous question and helps produce:

- A **decision frame**: what decision is being made, by whom, by when, against what criteria, with what constraints.
- An **issue map**: what must be true (or understood) to make the decision responsibly.
- A list of **alternatives**: what are the plausible options, including the status quo or phased approaches.
- A **tradeoff structure**: how options differ in benefits, risks, costs, dependencies, reversibility, and second-order effects.
- An **assumption register**: which premises are being relied upon, who owns them, and how they can be tested.
- A **verification plan**: what evidence is required to reduce uncertainty before the decision travels.

The essential point is that the model is not being asked to be "right." It is being asked to make the team's reasoning legible. A legible structure can be challenged. An illegible narrative cannot. That is why the scaffold matters: it externalizes cognition so that review becomes possible.

**The Structured Challenger.** The second abstraction is the model as a challenger—not in the sense of "debating for sport," but in the sense of systematic skepticism. Consulting teams are vulnerable to predictable distortions: confirmation bias, anchoring, availability, scope neglect, and premature convergence. A Reasoner can help counter those distortions by asking structured questions that humans do not reliably ask under pressure.

For example, the model can be instructed to:

- Identify the **weakest-link assumptions**: the small set of assumptions that, if wrong, would flip the decision.
- Surface **missing issues**: branches of the issue map that are absent or underdeveloped.
- Detect **assumption drift**: places where the output quietly introduced premises not present in the facts.
- Force explicit **status quo comparison**: how does "do nothing" perform relative to proposed action.
- Propose **disconfirming tests**: what evidence would falsify a favored option.

This challenger role is valuable because it makes skepticism a repeatable procedure rather than a personality trait. In many organizations, skepticism depends on one senior person asking the uncomfortable question. Level 2 aims to make uncomfortable questions part of the standard artifact. The model does not replace the senior person; it helps the senior person by ensuring the artifacts contain the prompts for skepticism.

**From narrative to inspectable objects.** The phrase "converts narrative into inspectable objects" is not just a metaphor. It is an operational instruction. In Level 1, the output is a document. In Level 2, the output is a set of objects that can be inspected independently. A tree can be checked for missing branches. A table can be checked for asymmetric treatment of alternatives. A register can be checked for ownership and testability. A verification plan can be checked for feasibility. These checks are hard to do when the logic is embedded in prose.

This is why Level 2 is so powerful in consulting work. It gives teams a way to create a structured record of their thinking that can survive scrutiny later. When a board asks "why did we think this was feasible?" the team can point to the assumption register and show what was believed, what evidence was sought, and what was verified at the time. That is a defensible posture. The model did not create defensibility; the model supported the creation of artifacts that enable defensibility.

**A Reasoner as a disciplined internal conversation.** Another useful way to understand Level 2 is to treat it as a structured internal conversation between roles. Even if you use a single model, you can prompt it to play distinct roles: strategist, skeptic, risk officer, operator, finance lens. The point is not to "multi-agent" the process (that is Level 3); the point is to ensure the reasoning bundle contains multiple lenses and is not just a single straight-line narrative. The model helps you stage the conversation in a consistent format: the strategist proposes a structure; the skeptic

challenges gaps; the risk lens identifies exposure; and the verification plan converts uncertainty into tasks.

This can be done within one prompt sequence, but Level 2 governance requires that the outputs remain explicitly provisional. Role-play does not create truth. It creates coverage of perspectives. Coverage is useful, but it can also be misleading if it is treated as validation. Therefore, the useful abstraction is always paired with the governance rule: perspective coverage is not evidence.

### 2.2.2 Dangerous misconception

The dangerous misconception at Level 2 is simple: *structured equals correct.* In consulting, structure is culturally associated with rigor. An issue tree feels like analysis because it resembles the artifacts that professionals use. But the model is not generating structure from verified data; it is generating structure from patterns in language. That distinction matters. The structure can be plausible, and plausibility is not correctness.

There are several specific misconceptions that follow from the "structured equals correct" error:

**Misconception 1: Exhaustiveness by formatting.** If an issue tree has clean branches, the team assumes it is MECE and exhaustive. But a tree can be beautifully formatted and still miss a critical driver—regulatory constraints, channel dynamics, supply chain fragility, stakeholder incentives, union considerations, or a capacity bottleneck. The tree is not proof of exhaustiveness; it is a hypothesis about what matters. If the tree is treated as exhaustive, the team can stop searching too early. That creates a subtle form of risk: not wrong answers, but *missing questions.* Missing questions are often the most expensive errors in strategy work.

**Misconception 2: Neutrality by symmetry.** A tradeoff table can give the appearance of neutrality because each option has pros and cons. But the model may describe one option's benefits in concrete terms and another's benefits in vague terms. It may overstate risks for options that are less common in its training distribution. It may produce "balanced" language that is actually biased by default. Symmetry of formatting is not symmetry of reasoning. Therefore, Level 2 requires a reviewer to check for asymmetric treatment: are we demanding more evidence for one alternative than another; are we granting feasibility to the option we already like.

**Misconception 3: Evidence by articulation.** The model can articulate a causal chain: "If we centralize procurement, we reduce unit costs through scale leverage." That is a plausible mechanism, but it is not evidence that it will occur in your organization. The model can describe a mechanism without knowing your contract structures, supplier fragmentation, implementation capability, or internal compliance constraints. At Level 2, causal articulation must be treated as a prompt for verification: "Is this mechanism present here?" If it is treated as evidence, the team will overestimate confidence.

**Misconception 4: Verification by planning.** A verification plan can be mistaken for verification. It is not. It is a list of tasks that *would* reduce uncertainty if executed. Teams sometimes feel reassured once the plan exists, even if it is never executed. That is a governance failure. Level 2 therefore requires ownership: verification tasks must have owners, evidence definitions,

and decision triggers. Without that, the plan is theater.

**Misconception 5: Decision legitimacy by documentation.** This is the decision laundering problem in a different form. A leader may say, "We have an issue map and tradeoff table, so we have done the analysis." The artifacts become a proxy for diligence. But diligence is not the presence of artifacts; diligence is the execution of verification and the quality of judgment. Level 2 artifacts are useful only insofar as they expose the logic and invite challenge. If they are used as stamps of legitimacy, they become actively harmful.

**Misconception 6: The model's confidence is meaningful.** Even with explicit instructions, models can produce language that implies confidence: "likely," "expected," "typically," "best practice." In consulting, those words often stand in for empirical claims. At Level 2, any such language must be treated as an assumption unless supported by evidence provided by the team. Otherwise, the model is smuggling in external "truth." A core governance move is to force these claims into the assumption register. If the model says "typically margins are X," that must become "assumption: typical margins are X; evidence needed: source; owner: analyst; status: unverified." Without this discipline, a Level 2 artifact becomes a vehicle for invented benchmarks.

**Why this misconception is particularly dangerous in consulting.** Consulting work is often reviewed by busy executives. Executives do not have time to read all the detail. They rely on signals: structure, clarity, professional tone. Level 2 outputs are rich in those signals. Therefore, if governance is weak, Level 2 can accelerate the propagation of error more than Level 1. A poorly written memo is sometimes ignored. A well-structured reasoning bundle may be adopted. That is why the mental model must be conservative: structure is a tool for scrutiny, not a substitute for it.

**The correction: treat structure as hypothesis.** The antidote to "structured equals correct" is to treat every reasoning object as a hypothesis artifact. The issue map is a hypothesis about what matters. The alternatives list is a hypothesis about the option set. The tradeoff mapping is a hypothesis about how options compare. The assumption register is a hypothesis inventory. The verification plan is a hypothesis-testing agenda. When teams adopt this posture, Level 2 becomes safe and powerful. When they do not, Level 2 becomes a factory of persuasive but unverified analysis.

### 2.2.3   Definition of a good Level 2 output

A good Level 2 output is not one that sounds impressive. It is one that makes the team's reasoning *inspectable*, makes uncertainty *explicit*, and makes verification *actionable*. The hallmark of a good Level 2 output is disciplined incompleteness: it does not pretend to know what it cannot know, and it refuses to cross the line into recommendation.

A good Level 2 output begins with a **clear decision frame**. Before the issue map, before the alternatives, the output must make explicit: what decision is being made, the decision maker, the time horizon, constraints (budget, timing, regulatory boundaries), and the criteria of success. Without a decision frame, reasoning artifacts float. They may be well structured but misaligned to the actual decision. Consulting often fails not because analysis is wrong, but because analysis

answers the wrong question. A decision frame is the simplest defense against that failure.

A good Level 2 output then provides an **issue map that is usable**. Usable does not mean perfect. It means scoped and prioritized. A model can generate a tree with dozens of branches, but a tree with too many branches is often a sign of low-quality structuring: it hides the signal in the noise. A good issue map separates what is core from what is peripheral, and it highlights where information is missing. It also labels assumptions explicitly when a branch depends on a premise. The map is not just a taxonomy; it is a working plan for inquiry.

Next, a good Level 2 output provides **alternatives that are meaningful**. That means alternatives that are distinct in mechanism or sequencing, not superficial variations. It includes the status quo when relevant, because many strategy decisions should be framed as "change versus not change," and failing to include the status quo is a common way to bias decision-making toward action. Alternatives should be written in a way that makes them testable: what would we actually do, what would change, what resources would be required.

Then, a good Level 2 output provides **tradeoffs that are asymmetry-aware**. The tradeoff section should not be a generic pros/cons list that would apply to any company. It should reflect the constraints and context in the facts provided. It should include second-order effects: what risks are created downstream, what dependencies emerge, what organizational friction might occur. It should also explicitly note where tradeoffs depend on assumptions. For example, "benefit: faster time-to-market" should be tied to assumptions about capability, vendor availability, or regulatory approval. Tradeoffs without explicit assumptions are rhetorical devices; tradeoffs with explicit assumptions are audit-ready reasoning.

The assumption register is the centerpiece. A good Level 2 output treats the **assumption register as a control mechanism**, not as a disclaimer. The register must be specific, owned, and testable. "Assume demand is strong" is not a good assumption. "Assumption: the target segment will adopt at a rate of X within Y months; owner role: commercial lead; test: pilot conversion data; evidence: cohort metrics; due date: placeholder" is closer to a good assumption. The purpose is not to be quantitative for its own sake, but to make assumptions operationally verifiable. A register also forces prioritization: which assumptions matter most. That leads to the weakest-link check: identify the assumptions that would flip the decision if wrong. This is what makes Level 2 valuable under time pressure: it tells you where to spend your verification time.

A good Level 2 output ends with a **verification plan that is executable**. It names tasks, owners, evidence definitions, and triggers. It should include what would change the team's mind. In consulting, "what would change our mind" is a powerful discipline because it resists confirmation bias. If a team cannot articulate what evidence would falsify an option, the team is not doing reasoning; it is doing advocacy. The verification plan does not have to be long; it has to be concrete.

Finally, a good Level 2 output is explicit about its posture. It labels itself *Verification status: Not verified. Human review required.*It avoids recommendation language. It avoids invented facts. It keeps the boundary between reasoning and decision intact. This is not bureaucratic. It is the difference between using AI as a disciplined thinking tool and using AI as a narrative machine that

produces decision-shaped artifacts.

The most important property of a good Level 2 output is that it *invites scrutiny.* It should make it easy for a reviewer to ask: "Where did this assumption come from?" "What evidence would test it?" "Are we missing an alternative?" "Did we treat options symmetrically?" "What is the weakest link?" If a Level 2 output makes those questions harder, it has failed, even if it looks professional. If it makes those questions easier, it has succeeded, even if it feels incomplete.

---

**Artifact (Save This)**

**Non-negotiable rules.**

1. **Facts are not assumptions.** Facts must be provided or verified; assumptions must be stated, owned, and testable.
2.
3.

**Required fields for Level 2 outputs (minimum).**

a. facts_provided (only what you supplied)

b. assumptions (explicit, owned)

c. issue_map (tree or bullets)

d. alternatives (2+ options, including status quo when relevant)

e. tradeoffs (pros/cons; risks; second-order effects)

f. open_questions (missing information)

g. verification_plan (what to check, how, who owns it)

h. verification_status: `"Not verified"`

---

## 2.3   What Level 2 can do, and what it cannot do

Level 2 is the first point in the maturity ladder where teams are tempted to treat the model as if it were doing "analysis." The temptation is understandable: the outputs look like the artifacts consultants use when they are actually thinking—issue trees, option sets, tradeoff matrices, assumption logs, and verification plans. But the governance-first posture of this book requires a disciplined separation between *reasoning support* and *decision ownership.* Level 2 is designed to improve the *structure* of thinking and the *auditability* of logic, not to substitute for judgment, evidence, or accountability. In practice, this section functions as an operating boundary: it tells the reader what to ask the Reasoner to do, what to forbid it from doing, and why those boundaries matter in consulting and corporate strategy work.

The core premise is simple: a Reasoner can be excellent at making thought visible, but it cannot know what is true unless you provide verifiable facts, and it cannot be responsible for decisions. Therefore, Level 2 is safe and useful when it produces artifacts that invite scrutiny and verification, and unsafe when it produces outputs that masquerade as conclusions. The difference is not merely

semantic; it determines how the output will travel through the organization. A properly governed Level 2 output slows down inappropriate certainty and accelerates the right kinds of work: clarifying the decision frame, surfacing missing information, prioritizing verification, and creating a structured basis for human deliberation.

### 2.3.1 CAN: reasoning tasks that benefit from explicit structure

Level 2 can do a great deal when the task is to impose structure on ambiguity. Consulting teams live in ambiguity: executives ask questions that are real but not fully specified, stakeholders disagree on objectives, and constraints emerge over time. The Reasoner is valuable precisely because it can help convert that ambiguity into inspectable objects. The goal is not to get an "answer," but to get a map of what would need to be answered, what could be chosen, what tradeoffs matter, and what assumptions must be tested.

**Decompose a question into an issue map and a decision frame.** One of the most valuable contributions of Level 2 is problem structuring. A decision frame makes explicit the elements that are often implicit: who is deciding, by when, against which criteria, under what constraints, and with what tolerance for risk. Once the frame exists, an issue map decomposes the problem into sub-questions that can be investigated and debated. This decomposition is not a bureaucratic artifact; it is the mechanism by which a team prevents scope drift and false completeness. In consulting terms, Level 2 can help the team move from "What should we do?" to "What must be true for any option to be viable?" and "What do we need to learn before we decide?" The model can propose a first-pass map quickly, which saves time, but the value comes from review: humans challenge the map, add missing branches, remove irrelevant ones, and confirm that the map matches the decision owner's criteria. The model is useful as a generator of candidate structure, not as the arbiter of correct structure.

**Enumerate plausible alternatives without selecting one.** Level 2 can produce a disciplined option set. This matters because many strategy failures begin with premature convergence: the organization falls in love with an option before it has defined what success means or what risks it is accepting. The Reasoner can help counteract this by forcing the team to articulate alternatives explicitly, including the status quo when relevant. A well-governed alternatives list does not merely restate the favored option in three disguises; it provides distinct mechanisms or sequencing choices: enter directly versus partner; centralize versus federate; phase pilots versus big-bang rollout; improve existing operations versus acquire capability. The model can also help articulate "hybrid" options, which are often neglected, and can help define disqualifiers and constraints that eliminate unrealistic options early. Importantly, Level 2 alternatives are framed as candidates for evaluation, not as recommendations. The correct output is "here are options and what would need to be true," not "here is what you should do."

**Map tradeoffs, constraints, and second-order effects.** Consulting decisions are rarely about a single metric. They involve tradeoffs: speed versus control, cost versus resilience, centralization versus autonomy, near-term earnings versus long-term positioning. Level 2 can help build a tradeoff

structure that makes these tensions explicit. The model can help list benefits and risks for each alternative, identify dependencies (what must be available, approved, or aligned), and surface second-order effects (how a decision changes incentives, capacity, or stakeholder behavior). This is particularly useful for strategy work because second-order effects are where many "obvious" solutions fail. A cost transformation may reduce spend but damage service levels. A centralization may create scale but slow decisions. A market entry may grow revenue but absorb management attention and create regulatory exposure. The Reasoner helps by forcing the team to articulate those effects as explicit items, which can then be tested or debated. The governance requirement is that these tradeoffs remain hypotheses unless supported by evidence. The model can propose plausible second-order effects, but humans must validate relevance in their context.

**Identify assumptions, weakest links, and information gaps.** This is the signature capability of Level 2. In many teams, assumptions remain implicit. They sit inside phrases like "we expect," "we believe," "likely," and "should be manageable." Level 2 can be used to extract those assumptions from the narrative and log them in an assumption register with ownership and testability. The model can also help prioritize assumptions by impact: which ones, if wrong, would flip the decision. This "weakest link" discipline is one of the most practical tools in consulting because it converts a complex, uncertain problem into a short list of high-value verification tasks. In addition, Level 2 can identify information gaps: what data is missing, what stakeholder views have not been collected, what operational constraints are unknown, and what external conditions (regulatory, market, competitive) remain unverified. Done well, this output gives the team a clear agenda: not "think harder," but "verify these three assumptions, ask these five stakeholders, and collect these two datasets."

**Produce a verification plan (what must be validated before a decision).** A verification plan is where Level 2 becomes governance, not just structure. The plan translates uncertainty into tasks with owners, evidence definitions, and timing. The Reasoner can generate candidate verification questions quickly: what evidence would confirm demand, what analysis would validate economics, what operational test would validate feasibility, what legal/compliance review is required, what pilot metrics would reduce uncertainty. But the plan must be owned by humans. A good verification plan assigns roles (not names in sensitive contexts), specifies what constitutes acceptable evidence, and defines triggers that would change the team's mind. The "change our mind" trigger is crucial: it prevents the verification plan from becoming performative. It forces the team to state upfront what evidence would cause them to stop, pause, or pivot. Level 2 can help draft this, but leadership must commit to it.

Across these CAN capabilities, the governance posture remains consistent: Level 2 creates structured reasoning artifacts that are **Not verified**. They are designed to be reviewed, challenged, and used to guide verification. They are not designed to be shipped as conclusions.

### 2.3.2   CAN'T: things you must not delegate to Level 2

The most common Level 2 failure is capability drift: a team begins by asking for structure, then slides into asking for answers. The slide is subtle because the outputs already look analytical. This subsection is therefore not an abstract warning; it is a hard operational boundary. If you cross it, you are no longer using Level 2 as a Reasoner. You are using it as an ungoverned decision engine, and you have lost defensibility.

**Choose an option, recommend a strategy, or provide "the answer."** Level 2 must not be allowed to select. The model can compare options, articulate conditions under which each option would be attractive, and highlight tradeoffs. But it must not recommend. In consulting, recommendations are not merely logical conclusions; they are professional judgments that integrate evidence, stakeholder context, risk appetite, and organizational constraints. A model does not hold accountability for the consequences of a recommendation. If a team asks the model to recommend anyway, the model will comply with a plausible narrative, and the organization may treat that narrative as an external "objective" voice. That is decision laundering. The governance-first rule is therefore strict: Level 2 may not pick winners. The most it may do is state, conditionally, "If these assumptions are true and these constraints hold, then option A would dominate on criteria X, but it carries risk Y." The decision remains human-owned.

**Invent facts, market benchmarks, competitor data, or citations.** This is non-negotiable. Strategy work is highly sensitive to fabricated data because small numbers can drive big decisions. Models can generate plausible benchmarks, pseudo-citations, and "industry norms" that sound credible. Even if the model is instructed not to invent, it can slip into phrases that imply external knowledge. The correct posture is to treat any such claim as an assumption unless the team provides the source and validates it. Level 2 is therefore prohibited from generating "market size," "typical margins," "competitor shares," or citations unless the inputs are supplied. In a governance-first workflow, the Reasoner can help formulate what data is needed, where it might be obtained, and how it would be used, but it must not fabricate it.

**Replace stakeholder judgment, political context, or client-specific constraints.** Consulting decisions occur inside organizations with incentives and history. A model can list "stakeholders" generically, but it cannot know who actually has power, where resistance will occur, what the CEO prioritizes this quarter, or what the board will tolerate. Those constraints are not details; they are the environment in which strategy is executed. If the Reasoner is asked to fill in these elements, it will produce generic narratives that may mislead the team. Worse, it can create a false sense that stakeholder alignment has been considered simply because it has been mentioned. Level 2 can support stakeholder mapping if the team provides the context, but it cannot replace judgment about what is politically feasible or ethically acceptable. Governance requires that stakeholder constraints be provided, owned, and reviewed by humans.

**Perform independent research unless you explicitly provide sources (and still requires verification).** Level 2, as defined in this book, assumes no tool use and no independent retrieval. If a team wants the model to incorporate external sources, that is already a step toward Level 3

workflows and must be governed accordingly. Even when sources are provided, the model must not be treated as a verifier. It can summarize or structure what the sources say, but a human must confirm the interpretation and the relevance. In consulting, source use is not only about accuracy; it is about appropriateness and applicability. A statistic from one market may not transfer. A "best practice" from one industry may fail in another. Therefore, even when sources exist, Level 2 outputs must remain *Verification status: Not verified. Human review required.*until a human has validated the evidence and its applicability.

**Additional practical boundaries (implied by the above).** To prevent drift, many teams adopt two additional prohibitions at Level 2. First, the model must not produce quantified forecasts or business cases unless the input data is supplied and clearly labeled as provided facts; otherwise, the numbers are too likely to be invented or misapplied. Second, the model must not produce language that implies certainty without a verification label. Words like "will," "definitely," or "proven" should be treated as red flags in Level 2 reasoning artifacts unless evidence is present and cited internally. These are not stylistic preferences; they are governance controls that protect the integrity of the consulting process.

In summary, Level 2 can do a great deal to improve consulting quality, but only if it is used for what it is: a structured reasoning assistant that produces inspectable artifacts and a verification agenda. It cannot do what organizations most want in moments of pressure: provide "the answer." The maturity ladder exists precisely to keep that pressure from turning a language model into an unaccountable decision-maker. Level 2 succeeds when it accelerates clarity and skepticism, and fails when it accelerates confidence without evidence.

## 2.4   The Level 2 governance wrapper (minimum controls)

Level 2 is where governance stops being an "AI hygiene" topic and becomes part of consulting professionalism itself. At Level 1, the primary control problem is straightforward: a chatbot can write convincingly, so the team must prevent invented facts, preserve authorship accountability, and ensure no draft leaves the room without review. At Level 2, the outputs look less like drafts and more like analysis: issue maps resemble problem structuring, tradeoff tables resemble diligence, and assumption registers resemble risk management. This is precisely why Level 2 needs an explicit governance wrapper. The danger is not only that the model may hallucinate a fact. The danger is that the model may produce *structured plausibility* that the organization mistakenly treats as verified reasoning.

The Level 2 wrapper is designed to do three things at once. First, it prevents **false rigor**: the phenomenon where structure creates unjustified confidence. Second, it prevents **assumption drift**: the gradual introduction of new premises that become accepted simply because they are written down. Third, it prevents **decision laundering**: the use of professional-looking artifacts to justify a choice that was made for other reasons. In consulting and corporate strategy, these are not hypothetical risks. They occur because time pressure, stakeholder dynamics, and executive

impatience reward outputs that look complete. Level 2 governance exists to ensure that what "looks complete" is still labeled and treated as *Verification status: Not verified. Human review required.*until evidence and judgment have been applied.

The wrapper is intentionally minimal. It is not a compliance program and it does not require new bureaucracy. It is a set of enforceable controls that can be applied by a consulting team on any engagement. The wrapper is also designed to be auditable: if the organization is later asked "what did we know, what did we assume, and what did we verify at the time," the wrapper produces artifacts that answer those questions. The wrapper turns Level 2 reasoning into a defensible workflow.

---

**Risk & Control Notes**

**Capability.** Structured reasoning artifacts (issue maps, alternatives, tradeoffs, assumption registers)

**Primary risks.** False rigor; hidden assumption drift; overconfidence from structure; premature convergence on a "best" option; scope creep; leakage of sensitive context

**Minimum controls.** Strict schema with no-recommendation rule; assumption register with ownership; weakest-link checks; verification planning; human review gate; versioning and prompt logs

---

The controls above can sound abstract until you translate them into operational rules that a team can actually follow. In Level 2, the single most important operational move is the **strict schema**. A schema is not a formatting preference; it is a control mechanism. When you force every output into fields such as facts_provided, assumptions, issue_map, alternatives, tradeoffs, open_questions, and verification_plan, you prevent narrative drift from hiding uncertainty. The schema forces the model to reveal whether it is using facts or creating assumptions, and it forces the team to confront what is missing. The schema also enables review: a reviewer can scan the assumptions and open questions quickly without reading every sentence of the narrative.

The second operational move is the **no-recommendation rule**. This rule must be enforced consistently, because the easiest way for Level 2 to go wrong is for the output to "pick a winner." Even when the model is asked not to recommend, it can imply a preference through language: describing one alternative as "pragmatic" and another as "risky," or using stronger verbs for one option. Governance must treat any preference cue as suspect. The goal of Level 2 is not to avoid judgment; the goal is to ensure judgment remains human-owned and evidence-informed. Teams should therefore treat the Reasoner's role as: structure, compare, stress-test, and propose verification. The moment the model begins to decide, the workflow becomes indefensible.

The third operational move is the **assumption register with ownership**. The assumption register is the "anti-hallucination" tool of Level 2, but it is also much more than that. It is the instrument that prevents implicit premises from becoming invisible drivers of a decision. In a typical consulting process, assumptions often appear as vague sentences: "We assume the market will grow,"

"We assume we can execute," "We assume stakeholders will align." Those assumptions are rarely tested because they are not operationally testable. A Level 2 assumption register forces specificity: what exactly is assumed, why it matters, what would falsify it, and who owns the verification task. Ownership matters because verification without ownership is theater. A register without owners is just a list of things that might be wrong. A register with owners becomes a workplan for reducing uncertainty.

The fourth operational move is the **weakest-link check**. This is a practical form of risk management tailored to strategy work. Instead of trying to verify everything, the team identifies the assumptions that would flip the decision if wrong. This is not just prioritization; it is a defense against false rigor. False rigor occurs when the team feels safe because many boxes have been filled. The weakest-link check keeps the team honest by asking: "If we could only verify three things, what would they be?" It also exposes where the decision is fragile. A fragile decision is not necessarily a bad decision, but it must be recognized as fragile so that leadership can choose a posture (pilot, stage-gates, reversibility, contingency plans) that matches the fragility.

The fifth operational move is **verification planning**. Level 2 ends with questions, not answers. The verification plan converts questions into tasks: what evidence to gather, which stakeholders to consult, which analyses to run, and what thresholds would change the team's mind. This is where Level 2 becomes defensible. A team that produces a brilliant issue map but no verification plan is still producing structured speculation. A team that produces a verification plan is producing a disciplined path to decision readiness.

The final two controls—**human review gates** and **versioning with prompt logs**—are what make the workflow auditable. Strategy work often lives in slides and spoken discussions. The reasoning that shaped a decision can be hard to reconstruct later, especially when outcomes disappoint. Level 2 governance insists that the reasoning artifacts and their evolution be saved as a "bundle" with version history and reviewer notes. This is not to create bureaucracy. It is to preserve accountability. If a decision is revisited, the organization should be able to say: "At the time, these were the facts provided, these were our assumptions, these were our open questions, and this was our plan to verify." That posture is what protects credibility.

### 2.4.1 The "Reasoning Bundle" (what you must save)

A Level 2 Reasoner is only as good as the record it leaves behind. In many consulting environments, the temptation is to treat reasoning as ephemeral: you do analysis, you present a recommendation, and you move on. But organizations increasingly demand traceability, especially when decisions are high-stakes or when outcomes deviate from expectations. The Reasoning Bundle is the minimum deliverable standard that makes Level 2 outputs reviewable, defensible, and reusable. It is also the mechanism by which a team can institutionalize "good Level 2 behavior" without relying on individual discipline alone.

The Reasoning Bundle has a simple philosophy: if the output influences decisions, the organization should be able to reconstruct how it was produced. This does not mean capturing every internal

conversation. It means capturing the minimal set of artifacts that show what the model was told, what it produced, what assumptions were introduced, what verification was planned, and what humans changed. The bundle functions like a lightweight audit trail for strategy reasoning.

---

**Artifact (Save This)**

**Minimum deliverable standard (Reasoning Bundle).**

1. Input snapshot (redacted): facts, constraints, objective, time horizon.
2. Prompt text: exact instruction used + any follow-ups.
3. Output JSON: full Level 2 schema + *Verification status: Not verified. Human review required.*
4. Assumption register: each assumption has an owner and a test.
5. Verification plan: tasks, owners, sources to consult, deadlines.
6. Reviewer notes: challenges, edits, and what was verified.
7. Version history: v0 (model) → v1 (human) → v2 (reviewed draft)

---

Each element of the bundle exists for a reason, and each is a control against a specific failure mode.

**1) Input snapshot (redacted).** This is the foundation. If the organization cannot see what facts and constraints were provided, it cannot evaluate whether the reasoning was appropriate. The snapshot must be redacted because consulting inputs often contain sensitive information, but it must still be sufficiently specific to preserve meaning. The purpose is not to store confidential data in the wrong place. The purpose is to establish a boundary: the model was allowed to use only these facts. If later the output contains claims not supported by the snapshot, the team can identify assumption drift immediately. The snapshot also enables peer review: another consultant can look at the inputs and ask whether the issue map and alternatives make sense given the stated objective and constraints.

**2) Prompt text.** The prompt is part of the professional work product. Two prompts that appear similar can produce meaningfully different outputs because the prompt defines the boundaries: whether the model is allowed to introduce assumptions, whether it must label uncertainty, and whether it must avoid recommendations. Capturing the prompt allows reproducibility. It also supports governance learning: teams can improve prompt patterns over time by reviewing what prompts produced good artifacts and which prompts produced risky outputs. In many firms, prompts become internal IP. The prompt log is the seed of that IP, and it must be versioned like any other method.

**3) Output JSON with *Verification status: Not verified. Human review required.*** The JSON output is the structured artifact. It is not enough to save the narrative memo; you need the discipline fields. The JSON structure ensures that facts, assumptions, open questions, and verification plans are not lost inside prose. It also ensures that outputs can be compared across iterations. For governance, the *Verification status: Not verified. Human review required.*label is not

optional. It is the mechanism that prevents the organization from treating the artifact as verified analysis. In practice, teams should treat the verification status as a gate: nothing moves forward without explicit status and reviewer sign-off.

**4) Assumption register with owners and tests.** The register is the center of gravity of Level 2 governance. It is the artifact that prevents the model's fluency from turning into unexamined premises. Ownership is critical. Without ownership, assumptions linger. With ownership, assumptions become tasks. A good register also includes testability: what evidence would confirm or falsify the assumption. Testability disciplines the team. It prevents vague assumptions from being used as scaffolding for confident conclusions. It also creates a natural bridge to Level 3, where verification tasks become workflow steps with checkpoints.

**5) Verification plan.** The plan operationalizes skepticism. It turns "we should validate" into "we will validate X by doing Y using evidence Z." This is where Level 2 outputs become aligned with consulting reality: time is limited, so verification must be prioritized. The plan should therefore reflect the weakest-link assumptions and the decision timeline. If the decision must be made in two weeks, the plan should not list tasks that would take three months. A verification plan that cannot be executed is worse than none, because it creates the illusion that due diligence exists. In governance-first practice, the plan should include owner roles, evidence definitions, and a due-date placeholder that forces the team to set a schedule.

**6) Reviewer notes.** Review is where accountability lives. If the model output is v0, reviewer notes are the record of professional judgment applied to that output. Reviewer notes should include: what was challenged, what assumptions were removed or revised, what additional facts were introduced (and their source), and what remains unverified. The notes do not need to be long; they need to be explicit. The purpose is to preserve the "why" of changes. Without reviewer notes, versioning is a sequence of edits without rationale, which is not defensible under scrutiny.

**7) Version history.** Versioning is the simplest way to prevent accidental loss of accountability. It preserves the original model output, the human-edited version, and the reviewed version. This allows the team to show how raw AI outputs were transformed into professional artifacts. It also supports internal learning: teams can compare v0 and v2 to see what kinds of model outputs consistently require correction. Over time, those patterns become training material and prompt improvements. Version history is also a control against "AI washing," where teams present a final output without acknowledging that it was heavily shaped by a model. In governance-first practice, transparency is part of trust.

Taken together, the Reasoning Bundle is a minimal system of record. It is small enough to be practical and strong enough to be defensible. It also creates a consistent habit across engagements: every Level 2 reasoning exercise produces the same categories of artifacts, making team review and leadership oversight easier.

### 2.4.2 Redaction and minimum-necessary input (Level 2)

If Level 1 is where teams learn not to paste confidential information casually, Level 2 is where teams must learn not to paste *strategically sensitive context* casually. The risk profile changes because Level 2 often tempts teams to share more: the problem feels complex, so they provide richer context to get better structuring. The governance-first posture is to treat this temptation as a control problem, not a personal failing. The right approach is to implement a "minimum-necessary input" discipline that is strong enough to protect confidentiality while still allowing the model to produce useful reasoning artifacts.

The principle is simple: provide only what is required to structure the reasoning, not what would be required to execute the strategy. In practice, that means redacting names, removing identifiers, and normalizing sensitive numbers. But at Level 2, it also means avoiding certain classes of information that are uniquely sensitive in strategy work: negotiation posture, pricing strategy, restructuring plans tied to individuals, internal politics, board dynamics, and unannounced capital allocation intentions. These are the kinds of details that, if leaked, can harm the organization even if no client names are present.

A governance-first Level 2 workflow therefore treats redaction as part of the method. It is not a separate compliance step. It is built into how the prompt is prepared. Teams should use placeholders (Client A, Competitor X, Region Y) and keep any mapping offline. Sensitive numbers should be normalized into ranges or indices when feasible. If an assumption depends on a specific number, the team can represent it as "Assume margin in the mid-teens" rather than pasting raw financial statements. The goal is to preserve the logic while reducing the confidentiality risk.

Level 2 also introduces a subtle risk: **identity inference**. Even if names are removed, a combination of details can re-identify a company or deal (industry, geography, transaction size, timeline, unique product features). Therefore, minimum-necessary input includes avoiding "unique fingerprints." If the reasoning does not require a specific detail, omit it. If the reasoning requires the structure of a contract, abstract it. If the reasoning requires a timeline, provide relative timing rather than exact dates, unless the dates are essential.

Redaction discipline also supports better reasoning. When teams abstract the problem, they are forced to clarify what actually matters. This reduces noise and helps the model focus on the decision frame and tradeoffs. A model can be distracted by irrelevant detail, just as humans can. Good redaction is therefore not only risk mitigation; it is analytical hygiene.

> **Checklist**
>
> **Level 2 redaction checklist (minimum).**
>
> 1. Replace client, competitor, and partner names with placeholders (Client A, Competitor X).
> 2. Remove deal terms, pricing, and negotiation positions unless approved.
> 3. Remove employee identifiers and sensitive HR/performance context.
> 4. Normalize sensitive numbers (ranges, indexed values) where feasible.
> 5. Keep a local mapping file (offline) if re-identification is required later.

In practice, teams should treat the checklist as a floor, not a ceiling. Two additional Level 2 redaction rules are often worth adopting as defaults. First, avoid pasting internal emails or raw transcripts. Summarize instead, because raw language often contains identifiers and incidental sensitive details. Second, avoid pasting content that reflects negotiation stance or board deliberations. Those materials are high-risk because their leakage can change counterpart behavior or create reputational harm. If such details are essential for reasoning, abstract them into constraints ("Assume counterpart is unwilling to accept price increase above X%") without including verbatim text or names.

Finally, redaction must be paired with environment discipline. A governance wrapper is weakened if the tool environment is uncontrolled. Level 2 reasoning should be performed only in approved environments consistent with organizational policy, with access controls and retention rules that match the sensitivity of the work. The chapter's focus is on method, not on vendor tooling, but the principle holds across tools: minimum-necessary input plus controlled environments plus saved bundles produces defensibility. Without these, Level 2 becomes a high-quality reasoning output attached to a low-quality governance posture, which is exactly the kind of mismatch that causes problems later.

The Level 2 governance wrapper is therefore not optional "risk language." It is the operational mechanism that allows consulting teams to use structured reasoning artifacts without accidentally converting them into unverified decisions. When applied consistently, it also becomes a professional advantage: teams can move quickly while maintaining a defensible record of what they knew, what they assumed, and what they verified. That is what governance-first maturity looks like at Level 2.

## 2.5 Core Level 2 workflow patterns (A–D)

Level 2 becomes practical only when it is routinized. The point is not to produce occasional "nice issue trees" that impress a partner in a meeting; the point is to create a repeatable discipline that any team can apply under time pressure without losing governance posture. The four workflow patterns below are intentionally simple and deliberately conservative. They are not a methodology that replaces consulting judgment. They are scaffolds that turn vague questions into structured artifacts, and structured artifacts into verification agendas.

Each pattern has the same governance backbone: (1) start with facts provided, (2) separate

assumptions, (3) keep outputs explicitly *Verification status: Not verified. Human review required.*, (4) refuse to recommend, and (5) convert uncertainty into a verification plan. Pattern A teaches the team to frame and decompose the problem. Pattern B forces exploration before convergence. Pattern C ensures that tradeoffs include second-order effects and failure modes rather than superficial pros/cons. Pattern D makes assumptions operational, identifies weak links, and produces the highest-value verification tasks. Used together, these patterns create a Level 2 "reasoning loop" that is defensible and auditable: frame → expand → compare → stress-test → verify.

### 2.5.1 Pattern A: Question-to-issue-map (problem structuring)

**Purpose.** Pattern A exists because most consulting questions arrive in the wrong form. They arrive as slogans ("Should we enter?"), impulses ("We need to cut costs"), or preferences ("We should centralize"), rather than as decision frames. Problem structuring is the foundational consulting skill: defining the decision, the criteria, the constraints, and the sub-questions that must be addressed. Level 2 can accelerate this skill by producing a first-pass decision frame and an issue map that the team then reviews and corrects. The goal is not to accept the model's tree as "the structure." The goal is to avoid starting from a blank page and to force explicit discussion of what matters before analysis begins.

**Inputs.** The inputs to Pattern A are intentionally limited. They include a messy question and any constraints the team already knows. If constraints are not provided, the model should surface them as open questions rather than inventing them. Typical inputs include:

- The raw question as asked by the client or executive sponsor.
- Objective and time horizon (even if tentative).
- Known constraints (budget, regulatory boundaries, capacity limits, timeline).
- Non-negotiables (e.g., "must maintain service level," "no layoffs," "must be reversible").

**Process.** Pattern A has three steps: frame, decompose, and expose gaps.

**Step 1: Frame the decision.** The decision frame is a short statement that fixes ambiguity. It answers: what decision is being made, who owns it, by when, against what criteria, under what constraints, and what constitutes an acceptable outcome. A good frame is specific enough to guide work but not so narrow that it embeds a solution. For example, "Decide whether to enter Market X in the next 12 months via build, partner, or acquire, given constraint Y and required return threshold Z." Even if return threshold Z is unknown, the frame can label it as an open question.

**Step 2: Decompose into an issue map.** The issue map translates the decision frame into sub-questions. A good Level 2 issue map has three properties. First, it is *hierarchical*: each branch breaks the problem into parts that can be addressed. Second, it is *decision-relevant*: branches map to the criteria and constraints in the decision frame. Third, it is *testable*: branches can be addressed through evidence, stakeholder input, or analysis, not through wishful thinking.

The model can generate a candidate map quickly, but the team must review it. Review should focus on missing branches and mis-scoped branches. In practice, the best review question is: "If

this branch is wrong or missing, could it change our decision?" If yes, fix it.

**Step 3: Surface open questions.** The most important output of Pattern A is often the list of open questions. These are not generic. They are the concrete items that prevent the team from pretending it knows more than it does. Open questions should be tied to the issue map and to the decision criteria. For example, "What is our actual execution capacity over the next two quarters?" is better than "What are the risks?"

**Outputs.** Pattern A produces three artifacts:

- A decision frame.
- An issue map (tree or structured bullets).
- A prioritized list of open questions and missing information.

All outputs must be explicitly *Verification status: Not verified. Human review required.*and must distinguish facts from assumptions.

**Governance notes.** Pattern A is where assumption drift begins if you are not careful. When the model sees an ambiguous question, it will try to complete it with typical constraints and typical criteria. The control is the schema: the model must place anything not provided into assumptions or open questions. Teams should also enforce a "no hidden objective" rule: the model must not assume what "success" means. If the client has not stated the criteria, the output must explicitly ask for them.

### 2.5.2   Pattern B: Alternatives-first (avoid premature convergence)

**Purpose.** Pattern B exists because consulting teams, like all humans, are vulnerable to premature convergence. The executive sponsor often has a preferred path. The team often has a favored framework. The organization often wants action. In that environment, analysis can become a rationalization exercise. Pattern B is a control against that. It forces the team to articulate alternatives before it invests heavily in validating one option. It also forces the team to include the status quo, which is the most common "forgotten alternative." Many decisions are framed as "we must do something," when in fact "do nothing" is a real option with real tradeoffs and often lower execution risk.

**Inputs.** Pattern B starts from Pattern A outputs:

- Decision frame.
- Issue map (at least the top-level branches).
- Constraints and non-negotiables.

If the team does not have a decision frame, Pattern B should not run. Without a decision frame, alternatives will be generic and may embed hidden criteria.

**Process.** Pattern B has four steps: generate, normalize, disqualify, and preserve optionality.

**Step 1: Generate alternatives.** The model can propose options quickly. The governance rule is that generation must not become recommendation. The output should list at least two meaningful alternatives plus the status quo when relevant. Alternatives should differ in mechanism,

sequencing, or risk posture, not just in adjectives. For example: "partner first then build," "pilot in one region," "acquire capability," "optimize current footprint." These are operationally distinct.

**Step 2: Normalize alternatives.** A common failure mode is that one alternative is described as detailed and actionable while others are vague. This creates bias. Pattern B therefore requires normalization: each alternative must be described using the same template, including (a) what would actually change, (b) required capabilities, (c) expected benefits (hypothesized), (d) key risks, and (e) dependencies. Normalization is not to add bureaucracy; it is to prevent narrative favoritism.

**Step 3: Explicit constraints and disqualifiers.** Alternatives must be tested against constraints. If a constraint makes an alternative impossible, that should be stated explicitly, not left as an implied weakness. Disqualifiers are useful because they prevent endless discussion of options that are politically or operationally infeasible. However, disqualifiers must be human-owned. The model can propose possible disqualifiers, but the team decides what actually disqualifies an option. If the disqualifier depends on an unverified claim, it must be logged as an assumption.

**Step 4: Preserve optionality.** In strategy work, committing too early is often the biggest cost. Pattern B encourages the team to preserve optionality by identifying phased approaches and reversible steps. The model can help propose a staged path: "pilot," "stage-gate," "prove economics," "scale." This is not a recommendation; it is an option structure that can reduce downside. The governance rule is to mark any such staging as conditional: "If we choose this path, staging could reduce risk by allowing verification before scale."

**Outputs.** Pattern B produces:

- A structured alternatives list (including status quo when relevant).
- For each alternative: a normalized description plus explicit constraints and disqualifiers.
- A shortlist of "real contenders" (if requested) *without selecting a winner*, typically framed as "options that remain viable given current constraints, pending verification."

**Governance notes.** Pattern B is where "decision laundering" often begins. A leader may say "list alternatives," but then treat the most flattering option as the answer. The control is twofold: the no-recommendation rule and symmetry review. Reviewers should scan for asymmetric language ("pragmatic," "best," "obvious") and require it to be rewritten into conditional, evidence-linked phrasing. The output should also explicitly state: "This is an option set, not a recommendation."

### 2.5.3 Pattern C: Tradeoff table + second-order effects

**Purpose.** Pattern C exists because most tradeoff analysis in practice is shallow. Teams list pros and cons, but they do not model second-order effects: how the decision changes incentives, capacity, risk exposure, and future options. In consulting, second-order effects are where strategy meets operations. They are also where elegant plans fail. A tradeoff table that ignores second-order effects is not merely incomplete; it can be actively misleading because it creates the impression that risks have been considered. Pattern C forces tradeoffs to include dependencies, reversibility, optionality, and "what breaks if we are wrong."

**Inputs.** Pattern C requires:

- A set of alternatives (from Pattern B).
- Decision criteria (explicit or as open questions).
- Constraints and key assumptions (at least first-pass).

**Process.** Pattern C has five steps: define criteria, map first-order tradeoffs, surface second-order effects, stress failure modes, and label evidence needs.

**Step 1: Define decision criteria.** Tradeoff tables are meaningless without criteria. If criteria are not provided, the output must treat them as open questions. Common criteria in consulting are growth, cost, risk, time-to-impact, resilience, stakeholder impact, and reversibility. But the correct criteria are the client's criteria. The model can propose candidate criteria, but the team must confirm them.

**Step 2: Map first-order tradeoffs.** For each alternative, map benefits and costs against criteria. The goal is clarity, not precision. If quantification is not supported by provided data, the output should remain qualitative and explicitly provisional. The model can help ensure that each criterion is addressed consistently for each alternative.

**Step 3: Surface second-order effects.** This is the core of Pattern C. Second-order effects include:

- Incentive shifts (who wins, who loses, how behavior changes).
- Capacity shifts (what new work is created, what load is removed).
- Governance shifts (decision rights, control points, accountability).
- Risk shifts (new exposures created, old exposures reduced).
- Option value (does the choice preserve flexibility or lock the organization in).

The model can propose plausible second-order effects, but governance requires that they be marked as hypotheses unless supported by evidence. The purpose is to create a checklist of what to examine, not to declare what will happen.

**Step 4: "What breaks if we are wrong?"** Pattern C forces an explicit stress question: if the key assumptions behind an alternative are wrong, what breaks operationally, financially, or reputationally? This is the fastest way to reveal fragility. It also helps teams design mitigation: pilots, stage-gates, contingency plans, and reversibility features. Again, mitigation proposals must be conditional and non-recommending: "Possible mitigations include..." with ownership and verification needs.

**Step 5: Evidence labeling.** A tradeoff table should explicitly indicate where evidence is missing. For example, "Benefit depends on assumption A; evidence needed: pilot metrics." This prevents the table from being read as final analysis. It also links Pattern C to Pattern D and to the verification plan.

**Outputs.** Pattern C produces:

- A tradeoff table (criteria by alternative).
- A second-order effects list per alternative.

- A failure-mode view: what breaks if key assumptions fail.
- A set of evidence needs tied to tradeoffs.

**Governance notes.** Pattern C is the pattern most likely to create false rigor. Tables look authoritative. Therefore the output must explicitly label itself as a hypothesis artifact, and it must avoid numeric precision unless numbers were provided. Reviewers should be trained to ask: "Which cells of this table are evidence-based and which are assumptions?" If the answer is unclear, the table must be revised.

### 2.5.4 Pattern D: Assumption register + weakest-link test

**Purpose.** Pattern D is the governance heart of Level 2. Every strategy decision rests on assumptions, but not all assumptions matter equally. The failure mode at lower maturity is that assumptions remain implicit, unowned, and untested. Pattern D makes assumptions explicit, assigns ownership, defines tests, and then performs a weakest-link check to prioritize verification. This pattern converts Level 2 from "structured thinking" into "structured skepticism," which is what makes the approach defensible.

**Inputs.** Pattern D draws from the prior patterns:

- Decision frame and issue map (Pattern A).
- Alternatives (Pattern B).
- Tradeoffs and evidence gaps (Pattern C).

**Process.** Pattern D has four steps: extract assumptions, structure the register, identify weakest links, and convert into verification tasks.

**Step 1: Extract assumptions.** Assumptions are often embedded in language. The model is useful here: it can scan the issue map, alternatives, and tradeoffs and identify implicit premises. For example: "We can hire talent," "Suppliers will comply," "Customers will switch," "IT can deliver," "regulators will approve," "stakeholders will align," "the organization can absorb change." The governance rule is that assumptions must be separated from facts. Any claim not present in facts_provided becomes an assumption, even if it sounds "reasonable."

**Step 2: Structure the assumption register.** A Level 2 register should include fields that force testability and ownership. At minimum:

- Assumption statement (specific).
- Why it matters (which decision branch depends on it).
- Owner role (who will validate).
- How to test (what evidence would confirm or falsify).
- Risk if wrong (what breaks).
- Status (unverified by default).

The act of structuring is itself a control: it prevents vague assumptions from hiding behind generalities.

**Step 3: Weakest-link test.** The weakest-link test answers: which assumptions are both

high-impact and highly uncertain? These are the assumptions that, if wrong, would flip the decision or make the preferred option infeasible. A useful way to do this is to rate assumptions on two axes (impact and uncertainty) and pick the top three. The model can propose candidates, but humans must confirm because impact depends on organizational context and decision criteria. The point is to prioritize verification. Teams should also use weakest-link results to discuss risk posture: if the decision depends on a single fragile assumption, leadership may choose staged commitment rather than full commitment.

**Step 4: Convert weakest links into verification tasks.** Pattern D ends with action. Each weakest-link assumption becomes a verification task with evidence definitions, owner roles, and timing. This is where Level 2 stops being "thinking" and becomes "work." Verification tasks may include stakeholder interviews, pilot experiments, financial analysis, operational capability reviews, legal/compliance consultations, or external research performed by humans using approved sources. The model can draft the plan, but humans own execution.

**Outputs.** Pattern D produces:

- A complete assumption register for the reasoning bundle.
- A list of weakest-link assumptions (top three by default).
- A verification task list derived from weakest links.
- An explicit statement of fragility: "This decision is sensitive to assumptions A, B, C."

**Governance notes.** Pattern D is the antidote to assumption drift and false confidence. However, it can fail if teams treat the register as a compliance artifact rather than a living tool. The control is to connect the register to decisions: reviewers should ask whether the verification plan is actually derived from the weakest-link assumptions. If not, the verification plan is likely performative. Another control is to ensure that assumptions do not quietly become facts in later drafts. Versioning helps: when a later draft uses an assumption as if it were true, the reviewer can catch the drift.

Taken together, Patterns A–D define the minimal Level 2 operating discipline. Pattern A frames and decomposes the decision. Pattern B forces exploration of options before convergence. Pattern C makes tradeoffs real by including second-order effects and failure modes. Pattern D turns assumptions into a prioritized verification agenda. None of these patterns requires the model to be "right." They require the model to be structured and honest about uncertainty. That is the governance-first promise of Level 2: not better answers, but better thinking records, better questions, and better decision hygiene.

## 2.6 Mini-cases (Level 2): structured reasoning in real consulting scenarios

The purpose of these mini-cases is not to demonstrate how a model can "solve" consulting problems. That would be both unrealistic and unsafe. The purpose is to demonstrate how Level 2 reasoning artifacts can be used to improve professional discipline: clearer decision frames, explicit assumptions,

symmetric alternatives, tradeoffs that include second-order effects, and verification plans that convert uncertainty into owned tasks. Each case below includes (1) a small, deliberately incomplete fact pattern, and (2) the expected shape of a Level 2 reasoning bundle. You should read these cases as templates. They are designed to be copied into your own work with redaction and minimum-necessary input, then used to generate *Verification status: Not verified. Human review required.*reasoning artifacts for human review.

A recurring governance theme runs through all four cases: Level 2 outputs are not recommendations. They are scaffolds for human judgment. The model must not pick an option, must not invent facts or benchmarks, and must not treat plausible-sounding logic as verified truth. The value is that the reasoning artifacts force a team to confront what it does not know, to surface the assumptions it is implicitly making, and to prioritize verification before commitments are made.

### 2.6.1 Case 1: Market entry (issue map + alternatives + verification plan)

**Small fact pattern (deliberately incomplete).** A North American industrial services company ("Client A") generates approximately \$1.2B in annual revenue, with 12% EBITDA margin, concentrated in two adjacent regions. Leadership is considering entering a third region ("Region X") where infrastructure spending is expected to increase over the next three years. Client A has no existing operations in Region X, but has relationships with national customers that operate there. The CEO wants an initial view within six weeks: whether to enter Region X and, if so, through what approach. Constraints: leadership does not want a major acquisition in the next 12 months; the organization is already running a large ERP modernization; and the CFO has stated that near-term cash flow must remain stable. There is limited internal data on Region X. The team has not yet defined target segments, pricing posture, or required return thresholds.

This fact pattern is intentionally under-specified. A Level 2 Reasoner should not "fill in" missing market data or competitor benchmarks. Instead, it should produce structured artifacts: a decision frame, an issue map, alternatives, tradeoffs, assumptions, open questions, and a verification plan.

**Level 2 objective.** Transform "Should we enter Region X?" into a decision-ready reasoning bundle that clarifies what must be learned, what options exist, and what verification is required.

**Expected reasoning artifacts (illustrative, not exhaustive).**

**Decision frame (example).** Decide whether Client A should enter Region X within the next 12–18 months, and if so through which entry mechanism (e.g., organic build, partnership, phased pilot), subject to constraints (no major acquisition in 12 months, limited change capacity due to ERP modernization, stable near-term cash flow). Define success criteria and decision owner timeline within six weeks.

**Issue map (top-level).** A Reasoner would typically propose branches such as:

- Market attractiveness (demand drivers, segment growth, customer willingness to pay).
- Right-to-win (capability fit, differentiators, brand/customer access).
- Economics (unit economics, cost-to-serve, pricing levers, cash flow impact).
- Operating feasibility (people, facilities, supply chain, service delivery model).

- Risk and governance (regulatory, safety, contract risk, reputational exposure).
- Entry path and sequencing (pilot vs scale, partnerships, investment staging).

Governance requires that any branch requiring unknown facts be reflected as open questions rather than asserted as "known."

**Alternatives (must include status quo when relevant).** A disciplined alternatives set might include:

a. **Status quo: do not enter Region X now.** Maintain focus on core regions; revisit after ERP stabilization.
b. **Phased pilot entry.** Start with a limited service offering for existing national customers in Region X; test demand and economics before scaling.
c. **Partnership-led entry.** Partner with a local operator to deliver services under a joint go-to-market; preserve cash flow and reduce operational load.
d. **Organic build-out (selective).** Establish a small footprint (sales + minimal operations) and scale as contracts are won; staged capex.

The Reasoner must not label any option "best." It should instead identify conditions under which each option is viable.

**Tradeoffs (example structure).** For each alternative, list benefits, risks, dependencies, reversibility, and second-order effects. For example:

- Pilot entry may reduce downside and preserve optionality, but may under-signal commitment to customers and risk slower scale.
- Partnership may accelerate access and reduce operational load, but introduces governance risk (quality control, margin sharing, brand risk).
- Organic build may preserve control and economics long term, but increases execution risk given ERP modernization and limited change capacity.

Second-order effects should be explicit: e.g., entering Region X could dilute management attention, create safety/compliance exposure, or change the firm's bargaining position with suppliers.

**Assumptions (examples).** A Level 2 output must separate assumptions from facts. Examples might include:

- Assumption: national customers will award incremental volume in Region X if Client A can service them.
- Assumption: service delivery can be staged without compromising safety and quality.
- Assumption: near-term cash flow impact can be managed via staged investment.

**Open questions (examples).**

- What are the prioritized target segments and service lines in Region X?
- What is the required return threshold and acceptable payback timeline?
- What is the current competitive landscape and customer switching cost?
- What is the feasible operating model given ERP program constraints?

**Verification plan (examples).** The plan should convert open questions into owned tasks:

- Customer validation: interviews with top 10 national customers; evidence: willingness-to-award, volume signals, contracting timelines.
- Economics validation: build unit economics for pilot service lines; evidence: cost-to-serve ranges, pricing bands from internal contracts.
- Feasibility validation: operations and safety review; evidence: staffing plan, compliance requirements, ramp timeline.
- Partnership scan: identify 3–5 local operators; evidence: capability fit, governance terms, reputational risk screening.

The output remains *Verification status: Not verified. Human review required.*until those tasks are executed and reviewed.

**Why this is Level 2 (not Level 1).** A Level 1 chatbot can draft a market entry memo. Level 2 produces the reasoning objects that make that memo defensible: it clarifies what must be true, what is assumed, and what must be verified before leadership commits.

### 2.6.2 Case 2: Cost transformation (tradeoffs + constraints + weakest-link test)

**Small fact pattern (deliberately incomplete).** Client A's board has asked for a cost reduction plan targeting \$80–\$120M run-rate savings over 18 months, with minimal impact to customer service and safety. The CEO believes "procurement and SG&A" hold the largest opportunity, while the COO warns that field operations are already stretched and that change fatigue is real. A prior cost program five years ago produced savings but damaged morale and created backlogs. Current constraints: unionized labor in two key geographies, a highly customized IT landscape, and a major customer renewal cycle in six months. Leadership wants an initial program hypothesis in four weeks, including a view on feasibility and risks. The team does not have detailed cost baselines yet, and the organization has not aligned on what "service levels" must be protected.

**Level 2 objective.** Produce a structured view of savings levers, tradeoffs, and feasibility constraints, and identify the weakest-link assumptions that could sink the program.

**Tradeoff structure (beyond generic pros/cons).** A Level 2 Reasoner can help structure levers into categories: procurement, SG&A, operations productivity, footprint optimization, and IT rationalization. But the governance posture requires that the Reasoner not invent savings amounts. Instead, it should produce:

- Lever hypotheses (what levers exist).
- Constraints and dependencies per lever.
- Second-order effects on service, safety, morale, and renewal cycle.

**Constraints and execution risk.** The core of this case is that cost transformation is rarely limited by "ideas." It is limited by execution capacity, sequencing, and organizational trust. A good Level 2 output therefore includes explicit constraints:

- Change capacity: concurrent initiatives, leadership bandwidth, frontline fatigue.

- Timing: customer renewals, union negotiations, seasonal demand.
- Governance: decision rights, approval pathways, and monitoring cadence.

**Alternatives framing (program archetypes, not "one plan").** Pattern B is useful here. The Reasoner can propose program archetypes:

a. **Fast-track procurement + SG&A** (low operational intrusion, quicker wins).
b. **Balanced program** (procurement + SG&A + selective ops productivity).
c. **Ops-led transformation** (deep operations redesign, higher savings potential, higher risk).
d. **Staged program** (phase 1 quick wins; phase 2 structural changes after renewal cycle).

The model must not recommend. It should instead articulate what must be true for each archetype to work.

**Weakest-link assumptions (examples).** This case benefits from Pattern D. Potential weakest links include:

- Assumption: procurement savings are achievable without service degradation (depends on supplier concentration and spec flexibility).
- Assumption: the organization can execute change without triggering morale collapse or attrition in critical roles.
- Assumption: unionized environments can accommodate productivity changes within the timeline.

The Reasoner should identify these as candidate weakest links, then convert them into verification tasks.

**Verification tasks (examples).**

- Build a baseline: confirm cost structure and spend categories; evidence: GL extracts, vendor lists, contract terms.
- Service-level definition: define "non-negotiable" service and safety metrics; evidence: SLA metrics, incident data, customer commitments.
- Capacity assessment: assess concurrent initiative load; evidence: PMO portfolio, critical path constraints.
- Labor constraints: union and HR review; evidence: contract provisions, negotiation timeline.

**Why this is Level 2.** The Reasoner is not producing a savings number. It is producing a structured program hypothesis, a risk-aware tradeoff map, and a prioritized verification plan that prevents a cost program from becoming a spreadsheet fantasy.

### 2.6.3 Case 3: Capital allocation (decision frame + scenarios + assumption register)

**Small fact pattern (deliberately incomplete).** Client A has $400M of projected free cash flow over the next two years. Leadership is debating three uses: (1) increase share repurchases, (2) invest in a new service line that could open growth in adjacent markets, and (3) reduce debt to improve resilience ahead of macro uncertainty. The CEO is optimistic about growth opportunities;

the CFO is concerned about volatility and wants a stronger balance sheet; the board has mixed views. Constraints: credit rating considerations, covenants, and a desire to maintain flexibility for potential M&A beyond 24 months. The team has incomplete forecasts and limited agreement on risk appetite. Leadership wants a board pre-read in three weeks that clarifies the decision and the tradeoffs, without committing to a specific recommendation yet.

**Level 2 objective.** Build a decision frame and scenario structure that clarifies how different capital allocation choices behave under uncertainty, and make assumptions explicit rather than smuggled into forecasts.

**Decision framing (must be explicit).** A good Level 2 output begins by stating the decision: "How should Client A allocate \$400M of projected free cash flow over two years across repurchases, growth investment, and debt reduction, given risk appetite, rating constraints, and desire for future optionality?" It should also specify criteria: shareholder value, resilience, flexibility, and strategic positioning. If criteria are contested, that must appear as an open question.

**Scenarios (structure, not numbers).** Level 2 can propose scenarios without inventing forecasts. For example:

- Base case: moderate demand growth, stable margins.
- Downside: demand shock, margin compression, refinancing constraints.
- Upside: strong adoption of new service line, improved pricing power.

The key is that scenarios must be described qualitatively unless the team provides quantitative inputs. The Reasoner can outline what variables drive the scenarios (volume, price, cost, working capital), but must not fabricate values.

**Tradeoffs under scenarios.** For each capital allocation option, the Reasoner can map:

- How the option performs in downside vs upside.
- What flexibility is gained or lost.
- What irreversibility exists (e.g., once repurchased, cash is gone).
- What signaling effects occur (market perception, board perception).

Second-order effects are important here: e.g., aggressive repurchases may reduce resilience, but may signal confidence; debt reduction may reduce return but improve strategic optionality.

**Assumption register (examples).** The capital allocation debate often hides assumptions about macro conditions, cost of capital, competitive dynamics, and growth feasibility. A Level 2 register might include:

- Assumption: the new service line can reach meaningful scale within 24 months.
- Assumption: maintaining the credit rating is a priority constraint (must be confirmed with board).
- Assumption: macro volatility warrants higher liquidity and lower leverage.

Each must have an owner role and a test (e.g., pilot evidence, board alignment, treasury analysis).

**Verification plan (examples).**

- Treasury analysis: covenant headroom and rating sensitivity; evidence: rating agency criteria,

stress tests.

- Growth diligence: service line economics and adoption proof; evidence: customer interviews, pilot results, competitive assessment.
- Board alignment: confirm risk appetite and priorities; evidence: structured board interviews, prior resolutions.

**Why this is Level 2.** The Reasoner helps build an honest decision conversation: it makes assumptions explicit, prevents numbers from being invented, and frames capital allocation as a set of tradeoffs under uncertainty rather than as a single "optimal" answer.

### 2.6.4 Case 4: Operating model redesign (options + governance implications)

**Small fact pattern (deliberately incomplete).** Client A has grown through acquisitions and now operates with fragmented functions across regions. Corporate leadership wants to redesign the operating model to improve consistency, reduce duplication, and strengthen controls. The COO favors regional autonomy because customer needs differ; the CFO wants more standardization in finance, procurement, and risk management; and the CHRO is concerned about talent retention and cultural backlash. Constraints: integration fatigue from recent acquisitions, uneven maturity of processes, and regulatory scrutiny in one business unit. Leadership wants a proposed operating model direction in eight weeks, but the organization has not agreed on decision rights, nor on how to measure "good governance" versus "bureaucracy."

**Level 2 objective.** Frame operating model choices as explicit options with governance implications, and identify the assumptions and verification needs that determine feasibility.

**Options (centralization vs decentralization is too simplistic).** Level 2 should avoid a false binary. A useful alternatives set might include:

a. **Federated model with standards.** Regions retain execution autonomy; corporate sets minimum standards, shared data definitions, and control points.
b. **Hybrid shared services.** Centralize transactional processes (finance ops, procurement ops) while leaving customer-facing decisions local.
c. **Functional excellence model.** Centralize key functions with strong business partnering; embed function leads in regions.
d. **Targeted centralization.** Centralize only high-risk or high-variance areas (e.g., compliance, safety, cybersecurity), leave others decentralized.

The Reasoner must not pick one. It should articulate conditions and tradeoffs.

**Governance implications (make them explicit).** Operating model is governance. The Reasoner can help surface:

- Decision rights: who decides what, and where approvals sit.
- Control design: what controls are required, how they are monitored.
- Accountability: what roles own outcomes, what KPIs align incentives.
- Risk posture: how regulatory scrutiny changes the required control set.

This is where second-order effects are critical: more centralization can reduce variance but increase cycle time; more autonomy can increase responsiveness but reduce control consistency.

**Assumptions and weakest links.** Operating model redesign often fails due to cultural and capability assumptions. Weakest links might include:

- Assumption: regions will accept corporate standards without active resistance.
- Assumption: shared services can deliver quality without losing local context.
- Assumption: data and process standardization is feasible given current IT fragmentation.

These should be converted into verification tasks: stakeholder interviews, process maturity assessment, IT architecture review, and pilot designs.

**Verification plan (examples).**

- Stakeholder alignment: structured interviews with regional leaders; evidence: willingness to adopt standards, identified red lines.
- Process maturity assessment: baseline variability and control gaps; evidence: process mapping, audit findings, KPI variance.
- IT feasibility: assess system constraints; evidence: architecture review, integration roadmap.
- Pilot governance: choose one function or region to pilot new decision rights; evidence: cycle time, control outcomes, adoption metrics.

**Why this is Level 2.** The Reasoner is used to make the operating model debate legible: explicit options, explicit governance tradeoffs, explicit assumptions, and an evidence plan to reduce uncertainty before redesign commitments are made.

Across all four cases, the same lesson repeats: Level 2 is not about producing recommendations. It is about producing disciplined reasoning artifacts that make recommendations possible later, after verification and human judgment. If you find yourself asking the model "which option should we choose," you have left Level 2. The correct Level 2 question is: "What must be true for each option to work, what are we assuming, and what do we need to verify before we decide?"

## 2.7 Risks and controls taxonomy (Level 2)

Level 2 is where governance stops being a set of polite disclaimers and becomes an operational necessity. The reason is straightforward: Level 2 artifacts look like analysis. An issue map resembles a consulting problem decomposition. An alternatives table resembles a structured evaluation. A tradeoff matrix resembles diligence. An assumption register resembles risk management. The visual and structural cues are strong enough that busy stakeholders may treat the package as "work completed," even when it is only "work structured." This section therefore provides a practical taxonomy: the predictable failure modes of Level 2, and the minimum control set that prevents those failure modes from turning into real organizational mistakes.

A taxonomy is not an academic exercise. It is a way of reducing ambiguity in supervision. When teams know what failure modes are typical, they can build review habits that detect them early.

When leaders know what controls are minimally required, they can enforce discipline consistently across engagements. The purpose is not to eliminate risk—consulting is inherently uncertain—but to ensure the organization does not confuse structured language with verified truth, or confuse a reasoning scaffold with a decision.

### 2.7.1 Failure modes you should expect

These failure modes are not rare. They emerge from the interaction of three forces: (1) the model's ability to produce plausible structure quickly, (2) human cognitive bias under time pressure, and (3) organizational incentives that reward closure. A disciplined team assumes these failure modes will occur and designs the workflow so that they are caught by default.

1. **False completeness: the issue map looks exhaustive but misses key drivers.**
   False completeness is the belief that because an issue map is neatly structured, it is complete. In practice, issue maps are always incomplete at first pass. The model generates what is typical, not what is decisive in your specific situation. It may omit constraints that are politically sensitive, operational bottlenecks that insiders take for granted, or external gating factors that are not salient in generic frameworks. The tree can be "balanced" while missing the one branch that would reverse the decision.

   False completeness often shows up in two patterns. First, the tree is too generic: it repeats a standard market-entry or cost-transformation template without anchoring to the constraints of the case. Second, the tree is too broad: it lists everything imaginable, creating the illusion of coverage while hiding the absence of prioritization. In both cases, the tree becomes a comforting artifact rather than a working plan.

   The governance response is to treat the issue map as a hypothesis and to require a "killer branch" review: what could kill the decision that is not represented? Teams should also require explicit linkage between decision criteria and tree branches. If a criterion has no branch, the map is incomplete. If a branch does not connect to a criterion, it may be noise or mis-scoped. Finally, teams should stress-test for context-specific drivers: regulatory gating, change capacity, timing constraints, stakeholder incentives, and irreversibility. These are the drivers most likely to be missing and most likely to matter.

2. **False rigor: tables and frameworks create unjustified confidence.**
   False rigor is the tendency to treat structured artifacts as if they were validated analysis. Tables are particularly dangerous because they visually encode comparability and evaluation. A tradeoff matrix can feel like the team has assessed options against criteria, even when the entries are unverified hypotheses. The model's fluent language amplifies this effect: if it can articulate a mechanism, humans often treat the mechanism as evidence.

   False rigor appears when the organization starts using Level 2 outputs as if they were final. For example, a tradeoff table is pasted into a board pre-read without labeling assumptions. Or an issue map is used to claim "we covered all areas" when the branches were never investigated. Or a verification plan exists on paper but was never executed, yet stakeholders feel reassured.

The governance response is explicit basis labeling. Every major claim in a Level 2 artifact must be classifiable as one of: (a) fact provided, (b) assumption, or (c) open question. If the team cannot classify a claim, it is not ready to travel. A practical control is to require that tradeoff tables include an "evidence needed" field or notation. Another control is to ban numerical specificity unless numbers were supplied and validated. Precision is a credibility amplifier; it must not be generated by the model.

3. **Assumption drift: the model introduces premises that become "accepted."**
   Assumption drift is when the model introduces new premises—explicitly or implicitly—and those premises become part of the team's worldview without being owned or tested. This is a Level 2 signature risk because the output is structured, and structure makes premises feel legitimate. Drift is rarely a dramatic hallucination. It is typically a plausible statement such as "customers will adopt quickly," "implementation can be staged," or "procurement savings are achievable." These statements are not obviously wrong. That is why they are dangerous: they slip through review and become the basis of downstream reasoning.

   Drift also occurs across iterations. In v0, the model states a premise as an assumption. In v1, a human edits for readability and moves it into narrative. In v2, someone treats the narrative as factual. Over time, the assumption has migrated into "what we know," even though no verification occurred. This is why versioning matters: it lets reviewers detect when assumptions are being silently upgraded.

   The governance response is the assumption register rule: every assumption must have an owner and a test. Ownership ensures accountability. Testability ensures the assumption can be verified or falsified. A good register is specific: it links the assumption to the decision branch it affects and specifies what evidence would settle it. A second control is "assumption quarantine": assumptions remain assumptions until evidence is produced. They do not become facts because they have been repeated.

4. **Premature convergence: one option becomes favored implicitly.**
   Premature convergence is when the team begins treating one option as the answer before the alternatives have been genuinely explored and before verification has reduced uncertainty. Level 2 can accelerate convergence because it produces artifacts that make one option feel more real. This happens when one alternative is described in richer detail, or when language subtly frames one option as "practical" and others as "risky." Even without explicit recommendation, the model can create a preference gradient.

   Premature convergence is often reinforced by organizational incentives. Executives want a direction. Teams want to show progress. A well-structured alternative set can become a stage for advocating the option that leadership already prefers. This can happen even unconsciously: people interpret the artifact as a justification tool rather than a reasoning tool.

   The governance response is alternatives-first discipline and symmetry enforcement. Require at least two meaningful alternatives plus the status quo when relevant. Require a normalized template for each option (what changes, dependencies, risks, reversibility, evidence needs).

Require explicit disqualifiers and treat any disqualifier as an assumption unless verified. A practical reviewer check is: "Did we make each option equally concrete?" If not, the artifact is biased by construction.

5. **Decision laundering: outputs are framed as "analysis" to justify a choice.**

   Decision laundering is the misuse of Level 2 artifacts to manufacture legitimacy. A decision is made for political, timing, or preference reasons, and the team then uses the model to create an issue map and tradeoff table that make the choice look evidence-based. The organization gets the comfort of "analysis" without doing verification. This failure mode is especially corrosive because it undermines learning. When outcomes disappoint, the organization cannot distinguish between a wrong assumption and an unexamined choice; it only has a narrative of diligence.

   AI makes decision laundering easier because it can produce professional-looking reasoning artifacts quickly, even when the underlying work was not done. The artifacts become proxies for diligence. This is not a problem that can be solved with better prompting alone. It requires governance: explicit labeling, reviewer skepticism, and a culture that treats verification as necessary for legitimacy.

   The governance response is threefold. First, enforce the no-recommendation rule so that the model cannot be used as an "objective voice" to select an option. Second, require verification planning and evidence gates: if a decision is presented as "analysis," the bundle must show what was verified. Third, require reviewer notes that document challenges to the preferred option and document what would change the team's mind. If the bundle cannot articulate disconfirming evidence, it is advocacy, not analysis.

6. **Confidentiality leakage: more context shared because it "helps reasoning."**

   Level 2 invites over-sharing. Teams want better structure, so they paste more context: internal emails, negotiation positions, pricing strategies, or sensitive organizational politics. The model may produce a better issue map, but the governance posture may be compromised. This is a classic capability-governance mismatch: better reasoning output achieved by increasing confidentiality risk.

   Leakage can be direct (pasting names, deal terms, identifiable numbers) or indirect (providing a set of unique details that makes re-identification possible even without names). It can also be procedural: using unapproved environments or failing to apply minimum-necessary input discipline.

   The governance response is to treat redaction as part of the method, not as a separate compliance step. Use placeholders. Normalize numbers. Remove negotiation stance unless explicitly approved. Avoid raw transcripts; summarize. Maintain an offline mapping file if re-identification is needed. And ensure Level 2 work occurs only in approved environments consistent with firm policy. If the team cannot operate in an approved environment, it should reduce inputs, reduce sensitivity, or not use the tool.

### 2.7.2 Minimum control set (practical and enforceable)

The minimum control set below is designed to be enforceable in real consulting conditions: short timelines, messy inputs, and multiple stakeholders. These controls are not "nice to have." They are the boundary conditions for safe Level 2 usage. If a team cannot implement them, it should not treat Level 2 outputs as decision-support artifacts.

The controls can be understood as a chain. The schema creates structure. The no-recommendation rule preserves human decision ownership. The assumption register and weakest-link checks prevent implicit premises from driving choices. The human review gate ensures accountability. Versioning preserves traceability. Confidentiality guardrails protect sensitive information. If any link is missing, the chain fails.

---

**Checklist**

**Minimum controls (use in every engagement).**

1. **Schema requirement:** facts_provided / assumptions / issue_map / alternatives / tradeoffs / open_questions / verification_plan / verification_status.
2. **No recommendation rule:** the model must not select an option.
3. **Assumption register rule:** every assumption has an owner + test.
4. **Weakest-link check:** identify the top assumptions that could flip the decision.
5. **Human review gate:** decisions require human sign-off; AI outputs are supporting artifacts only.
6. **Versioning:** preserve v0 and track changes; keep prompt logs.
7. **Confidentiality guardrails:** minimum-necessary input; redaction; approved environments only.

---

Each control is simple on paper, but it matters because it blocks a specific failure mode.

**1) Schema requirement.** This control blocks false rigor and assumption drift. When the output must declare facts, assumptions, open questions, and a verification plan, it becomes harder for the model to smuggle in premises and harder for humans to treat the artifact as complete. The schema also enables quick review. A reviewer can scan assumptions and open questions first, rather than being seduced by narrative flow.

**2) No recommendation rule.** This control blocks premature convergence and decision laundering. The model cannot be allowed to "pick." Even a subtle preference cue can be misused as legitimacy. The team can still reach a recommendation, but it must be owned by humans and supported by evidence. The model's role is to structure reasoning, not to decide.

**3) Assumption register rule (owner + test).** This control blocks assumption drift. It forces assumptions to be named, owned, and testable. Ownership ensures someone is responsible for verification. Testability ensures assumptions are not vague disclaimers. If an assumption cannot be tested, it must be treated as a risk that shapes decision posture (e.g., staging, reversibility), not as a hidden premise.

**4) Weakest-link check.** This control blocks false completeness by prioritizing what matters. It forces the team to identify which assumptions are both high-impact and high-uncertainty. Those assumptions become the focus of verification. The weakest-link check prevents teams from doing "a little bit of everything" and calling it diligence. It makes verification efficient.

**5) Human review gate.** This is the accountability control. Level 2 artifacts are *Verification status: Not verified. Human review required.*by default. A human reviewer must challenge the artifacts, verify key claims, and decide what can travel. The gate also ensures that recommendations and commitments remain human-owned. In practice, the gate should include documented reviewer notes: what was challenged, what changed, and what remains unverified.

**6) Versioning and prompt logs.** This control makes the process reconstructable and prevents silent upgrades of assumptions. Preserve v0 model outputs. Track changes across iterations. Keep the exact prompts and follow-ups. This enables auditability and improves learning. It also protects the team when decisions are revisited: you can show what was known and what was assumed at the time.

**7) Confidentiality guardrails.** This control prevents a capability-governance mismatch. Better reasoning is not worth improper data exposure. Minimum-necessary input and redaction are required at Level 2 because strategic context is sensitive. Approved environments reduce risk of unintended retention or access. If guardrails cannot be met, reduce inputs, abstract the problem, or do not use the tool.

A final point: controls only work when they are practiced. The most effective way to institutionalize Level 2 discipline is to make the Reasoning Bundle the default work product and to make the checklist above the default review rubric. When every engagement produces the same structured artifacts, review becomes easier and governance becomes habitual. Level 2 maturity is not achieved by having one expert who "knows how to prompt." It is achieved by making disciplined reasoning and explicit verification a routine part of consulting work.

## 2.8   Prompt patterns and exercises (copy/paste)

Level 2 only becomes "real" in an organization when it becomes repeatable. The promise of a Reasoner is not that it will occasionally generate a clever issue tree. The promise is that it can be used as a consistent discipline: every time a team faces an ambiguous question, it produces the same class of reasoning artifacts, with the same governance posture, so that review becomes easier and decision hygiene improves over time. Prompt patterns are the mechanism by which this discipline is transmitted across teams. They are not magic words. They are standardized operating procedures: a small set of reusable prompt templates that (1) force facts/assumptions separation, (2) forbid recommendation and invention, (3) produce inspectable objects, and (4) end with a verification plan rather than a conclusion.

This section provides three copy/paste prompts plus a training exercise set. The prompts are designed to be used in a governed environment and to produce a Reasoning Bundle consistent with

the minimum control set in this chapter. A crucial note: the templates below include the word "STRICT," but strictness is not a wish; it is a constraint the model may fail to follow. In practice, teams should treat these prompts as part of a larger wrapper that validates JSON and enforces the schema. Even when you are working manually, you should enforce a "schema compliance check" before using the output: if the model produces narrative rather than structured fields, the correct response is to re-run the prompt with a stronger constraint or to manually restructure the output. Do not accept a half-structured output as "good enough," because half-structured outputs are how assumptions and invented facts hide.

A second note: these prompt templates assume Level 2 posture—no tool use, no browsing, no retrieval. If you want the model to use external sources, that moves you into Level 3 workflows and requires additional governance (tool logs, source capture, routing, and checkpointing). The prompts below are therefore intentionally self-contained and conservative: they force the model to treat anything not provided as unknown.

## 2.8.1 Prompt Template 1: Governance-first reasoning wrapper (universal Level 2)

Template 1 is the universal wrapper for Level 2. It is what you use when you want to convert an ambiguous question into a structured reasoning bundle: decision frame, issue map, alternatives, tradeoffs, assumption register, weakest-link assumptions, open questions, and a verification plan. It is also the template that most strongly enforces the non-negotiable rules: do not recommend, do not invent, be explicit about uncertainty, and separate facts from assumptions.

There are three practical implementation notes that matter for this template.

First, keep the context short and structured. The model performs best when you provide facts as bullets and constraints explicitly. Avoid pasting raw emails or transcripts. Summarize instead. This is not only a confidentiality control; it is a quality control. Raw text often contains irrelevant detail that distracts the reasoning.

Second, treat "owner$_r$ole¡asarole, notaperson.Level2iswhereteamsaretemptedtopastenames.Resistthat.Userole$Partner, EngagementManager, CFO, HeadofOps, LegalCounsel, PMOLead.Thispreservesconfidentialitywhi$

Third, keep "due$_date_p$laceholder¡asaplaceholderunlessyouarecertainandcomfortable.Inaconsultingsetting, d$

**Prompt (Copy/Paste)**

```
ROLE: You are a reasoning assistant for management consulting deliverables.
Do not recommend. Do not invent facts, data, benchmarks, or citations.
Be explicit about uncertainty. Separate facts from assumptions.

CONTEXT (facts provided by user; redacted):
[PASTE FACTS + CONSTRAINTS + OBJECTIVE HERE]

TASK:
1) Build an issue_map for the question below.
2) List alternatives (2+), including status quo when relevant.
3) Provide tradeoffs and key risks for each alternative.
4) Create an assumption register and identify the weakest-link assumptions.
5) Produce a verification_plan with concrete questions to verify (no guessing).

QUESTION:
[PASTE THE CLIENT/TEAM QUESTION HERE]

OUTPUT FORMAT (STRICT JSON):
facts_provided: [bullet list]
assumptions: [bullet list]
issue_map: [tree or nested bullets]
alternatives: [list]
tradeoffs: [per alternative: pros/cons/risks/dependencies]
assumption_register: [assumption | why_it_matters | owner_role | how_to_test]
weakest_link_assumptions: [top 3]
open_questions: [list]
verification_plan: [task | owner_role | evidence_needed | due_date_placeholder]
verification_status: "Not verified"
questions_to_verify: [bullet list]
```

**How to review Template 1 outputs (practical rubric).** A team that uses Template 1 well should adopt a fast review routine. The reviewer should not start by reading the issue map. They should start by scanning for governance red flags:

1. Do facts_provided contain only what the user supplied, or did the model "add context"?
2. Are assumptions explicit and do they include premises that look like facts or benchmarks?
3. Are alternatives symmetric and meaningful, including status quo when relevant?
4. Are tradeoffs conditional and hypothesis-labeled, or do they read like conclusions?
5. Does the assumption register include owner roles and test methods, or is it vague?
6. Are weakest-link assumptions truly decision-flipping, or are they generic?
7. Does the verification plan contain concrete evidence requests, or does it say "research more"?

If any of these checks fail, the correct action is to revise the prompt or manually correct the output,

and then preserve the corrected version as v1 with reviewer notes. Do not let a flawed v0 artifact travel as if it were disciplined.

### 2.8.2 Prompt Template 2: Stress-test the reasoning (red-team pass)

Template 2 is the "skeptical reviewer" pass. It is the simplest way to institutionalize professional skepticism without needing a separate person to always play the contrarian. The aim is not to defeat the reasoning; the aim is to identify where the reasoning is fragile, incomplete, or biased by implicit premises. In other words, Template 2 is your systematic defense against false completeness, false rigor, and premature convergence.

There are four implementation notes.

First, Template 2 must not add facts. It is a critique pass, not an enrichment pass. When reviewers allow the model to "fix" gaps by inventing facts, they create a loop where the model hides its own uncertainties.

Second, the input should be the structured artifacts, not the narrative. If you paste narrative prose, the critique will be less targeted and more rhetorical. If you paste structured fields, the critique becomes an inspection of objects.

Third, the output is deliberately short. Its job is to produce an actionable list of gaps and high-value verification tasks. It should not produce a new full reasoning bundle; that would blur responsibilities. If you want a revised bundle, you run Template 1 again with the gaps addressed.

Fourth, teams should treat this as mandatory for any high-stakes decision or any time there is strong executive preference for one option. Bias and premature convergence are more likely under those conditions.

---

**Prompt (Copy/Paste)**

```
ROLE: You are a skeptical reviewer. Do not recommend. Do not add new facts.

INPUTS:
- facts_provided
- assumptions
- issue_map
- alternatives
- tradeoffs

TASK:
Identify:
1) Missing issues / non-MECE gaps
2) Hidden assumptions (unstated premises)
3) Where tradeoffs are asymmetric or misleading
4) What evidence would most reduce uncertainty

OUTPUT FORMAT (STRICT):
missing_issues:
hidden_assumptions:
tradeoff_gaps:
highest_value_verifications:
verification_status: "Not verified"
```

---

**How to use Template 2 in practice.** A disciplined team uses Template 2 as a gate: no Level 2 reasoning artifact is allowed to influence a recommendation memo or a board deck until it has been stress-tested. The stress-test output becomes part of the Reasoning Bundle as reviewer notes or as an explicit "red-team pass." The team then updates the issue map and verification plan accordingly and logs the change. Over time, this creates organizational learning: teams become better at anticipating gaps and writing stronger initial prompts.

### 2.8.3 Prompt Template 3: Convert reasoning into an executive discussion guide

Level 2 artifacts can be too technical or too "internal" to travel directly into executive discussion. Executives do not want a full issue map. They want the decision questions, the uncertainties, the evidence requests, and the risks that need explicit acknowledgment. Template 3 converts structured reasoning artifacts into a 30-minute discussion guide. The output is still not a recommendation. It is a meeting tool: it helps a team use a short executive session to clarify criteria, surface constraints, and assign verification requests.

Template 3 matters because good consulting is not merely analysis; it is guided decision conversation. Many engagements fail because the team does analysis in isolation and then presents conclusions without first aligning on criteria and risk appetite. A discussion guide can prevent

this. It forces the team to ask decision owners what matters before committing to a direction. It also helps create a record: the output can be attached to meeting notes and used as part of the verification plan.

Implementation notes:

- Keep the meeting guide neutral. It must not nudge the executive toward an option.
- Use questions that reveal priorities. "What matters most?" is not a good question. "Which criterion dominates if tradeoffs emerge: speed, resilience, or margin?" is better.
- Make verification requests concrete and owned. Executives are more likely to authorize verification when it is phrased as a specific evidence request.
- Include "what would change our mind" triggers. This is the strongest anti-laundering control because it commits the group to evidence that could falsify a preferred option.

---

**Prompt (Copy/Paste)**

```
ROLE: You are a consulting communication assistant.
Do not recommend and do not add facts.

INPUT:
[PASTE Level 2 reasoning artifacts]

TASK:
Create an executive discussion guide for a 30-minute meeting:
- 3 decision questions
- 5 clarification questions
- 5 verification requests (evidence to request)
- 3 risks to flag explicitly
- 3 "what would change our mind" triggers

OUTPUT FORMAT:
decision_questions:
clarification_questions:
verification_requests:
risks_to_flag:
mind_change_triggers:
verification_status: "Not verified"
```

---

**How to use Template 3 well.** Teams should treat the discussion guide as a pre-read for a checkpoint meeting, not as a substitute for analysis. The guide works best when paired with a short "facts and constraints" summary so that executives can react to what is known versus unknown. After the meeting, the team should update the decision frame, assumptions, and verification plan based on what was clarified. This is a key Level 2 discipline: executives often change the problem definition, and the artifacts must be revised accordingly. The updated artifacts should be saved as a new version with meeting notes attached.

### 2.8.4 Exercise set (team training)

Level 2 capability is not acquired by reading about it. It is acquired by repeating a small set of disciplined moves until they become muscle memory. The exercises below are designed for team training. They are tied to the recurring mini-cases in this chapter (market entry, cost transformation, capital allocation, operating model redesign). Each exercise is structured to produce a Reasoning Bundle and to force a governance review. In other words, the exercise is not "get a good answer." The exercise is "produce an auditable reasoning process."

A recommended training format is a 60–90 minute session with three roles:

- **Driver:** runs the prompts and assembles the Reasoning Bundle.
- **Reviewer:** enforces the minimum control set and writes reviewer notes.
- **Sponsor:** plays the executive and answers clarification questions.

Rotate roles across exercises. The goal is to make every team member competent in both production and review.

**Exercise 1: Question-to-issue-map under ambiguity (Pattern A).** Use Case 1 (market entry). Provide only the fact pattern and constraints. Run Template 1 to produce a decision frame and issue map. The reviewer must identify at least three missing issues that could flip the decision. Then re-run Template 1 with the missing issues incorporated. Save v0 and v1 with reviewer notes. Success criterion: the final issue map clearly distinguishes facts, assumptions, and open questions, and does not smuggle in market data.

**Exercise 2: Alternatives normalization and symmetry check (Pattern B).** Use the same market entry case. Generate alternatives including status quo. The reviewer must enforce normalization: each alternative must use the same template. The reviewer must also identify any language that implies recommendation and require revision. Success criterion: options are equally concrete and none is implicitly favored.

**Exercise 3: Tradeoff table evidence labeling (Pattern C).** Use Case 2 (cost transformation). Produce a tradeoff table comparing at least three program archetypes. For each tradeoff claim, label whether it is a fact, an assumption, or an open question. Where evidence is needed, specify the evidence type. Success criterion: the team can explain what would have to be verified for the table to become decision-support rather than speculation.

**Exercise 4: Assumption register and weakest-link test (Pattern D).** Use Case 3 (capital allocation). Build an assumption register with owner roles and how-to-test fields. Identify the top three weakest-link assumptions. Convert them into verification tasks with evidence needed. Success criterion: the weakest links are genuinely decision-flipping and the verification tasks are concrete enough to execute.

**Exercise 5: Red-team pass and revision loop (Template 2).** Take any Reasoning Bundle from Exercises 1–4. Run Template 2. The reviewer must decide which critique items require immediate revision versus which become future risks. Update the artifacts accordingly and record changes. Success criterion: the team demonstrates an explicit loop from critique to revised reasoning

artifacts, with version history intact.

**Exercise 6: Executive discussion simulation (Template 3).** Using a revised Reasoning Bundle, run Template 3 to create a 30-minute discussion guide. Have the sponsor role answer the questions as an executive would (including ambiguity and preference). Update the decision frame and verification plan based on that "meeting." Success criterion: the team updates artifacts based on sponsor input without turning the output into a recommendation.

**What to save for every exercise (Reasoning Bundle discipline).** Every exercise must produce a saved bundle:

- Input snapshot (redacted)
- Prompt text and follow-ups
- Output JSON (v0)
- Red-team output (if used)
- Revised output (v1) with reviewer notes
- Verification plan with owners and evidence needed

If a team completes these exercises and can produce consistent bundles, it has achieved practical Level 2 maturity. The organization will not merely be "using AI." It will be using AI in a way that makes thinking inspectable, reduces persuasive error, and strengthens the defensibility of strategy work.

## 2.9 Conclusion and transition to Level 3 (Agents)

Level 2 is the point in the maturity ladder where AI becomes genuinely useful for consulting work, and also the point where it becomes genuinely dangerous if the team mistakes structure for truth. The defining achievement of Level 2 is that it makes reasoning inspectable. Instead of relying on persuasive narrative, teams produce objects that can be reviewed: a decision frame that clarifies what is being decided, an issue map that clarifies what must be understood, alternatives that prevent premature convergence, tradeoffs that expose constraints and second-order effects, and an assumption register that forces hidden premises into the open. When used correctly, Level 2 does not make teams "smarter" in the abstract. It makes them more disciplined. It makes their thinking more legible to one another, more reviewable by leaders, and more defensible when decisions are revisited later.

But Level 2 also reveals a governance paradox: the better the artifact looks, the easier it is to misuse it. A clean issue map can create false completeness. A tradeoff table can create false rigor. An assumption register can create the illusion of risk management even when no one owns verification. And because these artifacts resemble the outputs of professional analysis, they can be used—intentionally or unintentionally—to launder decisions that were made for other reasons. In other words, Level 2 can turn "persuasive prose" into "persuasive structure." That is why the governance-first posture remains non-negotiable at this level: *Verification status: Not verified. Human review required.* is not a disclaimer; it is the boundary between reasoning support and

decision ownership.

### 2.9.1 Summary of main takeaways

The first takeaway is that Level 2 changes the unit of work. At Level 1, the unit of work is a draft: an email, a memo, a storyline. At Level 2, the unit of work is a reasoning bundle: decision frame, issue map, alternatives, tradeoffs, assumptions, and verification plan. This matters because it changes how consulting teams review and communicate. A manager can review a draft for clarity. A manager must review a reasoning bundle for logic, symmetry, missing issues, and untested premises. Level 2 therefore requires a stronger review habit. You cannot govern Level 2 by skimming for tone; you govern it by inspecting assumptions and evidence needs.

The second takeaway is that Level 2 is most valuable when it is used to *surface uncertainty* rather than to hide it. Good Level 2 outputs are incomplete by design. They expose open questions and they convert those open questions into verification tasks. In consulting, the most costly errors are often not "wrong answers," but unasked questions. Level 2's power is that it makes unasked questions visible and gives the team a structured agenda for reducing uncertainty before commitments are made.

The third takeaway is that the assumption register is the governance heart of Level 2. Strategy work always relies on assumptions. What changes at Level 2 is that assumptions become explicit, owned, and testable. This is what prevents the model from smuggling in premises and what prevents the team from accidentally treating plausible-sounding statements as facts. The weakest-link test operationalizes this discipline further: it forces prioritization. Instead of trying to verify everything, the team identifies the small set of assumptions that would flip the decision if wrong and focuses verification there. This is how Level 2 stays practical under time pressure.

The fourth takeaway is that Level 2 requires a strict boundary against recommendation. Many organizations want AI to "tell us what to do." Level 2 explicitly resists that impulse. The Reasoner can compare options, articulate conditions under which each option is viable, and identify risks and dependencies. But it must not select. Selection is not merely a logical operation; it is a professional judgment that integrates evidence, stakeholder context, and accountability. The no-recommendation rule is therefore not a limitation of the tool; it is a protection of professional responsibility.

The fifth takeaway is that process documentation is not optional if you want defensibility. Level 2 artifacts must be saved with prompt logs, reviewer notes, and version history. This is the minimal system of record that prevents assumption drift across iterations and makes later review possible. Without this record, Level 2 becomes a set of persuasive documents floating without provenance. With this record, Level 2 becomes a disciplined workflow that an organization can trust and improve over time.

Finally, the central risk at Level 2 is not that the model occasionally says something wrong. The central risk is that the organization begins to treat AI-generated structure as a substitute for verification and judgment. False rigor and decision laundering are the signature failure modes. The antidote is the minimum control set you have practiced throughout this chapter: strict schema,

assumption ownership, weakest-link checks, verification planning, and human review gates. If those controls are consistently applied, Level 2 becomes a professional advantage rather than a governance liability.

### 2.9.2 What comes next (preview of Level 3)

Level 3 is not "more intelligence." It is more process. If Level 2 is about making reasoning inspectable, Level 3 is about making the workflow inspectable. The capability shift is that the organization stops treating a single model call as the unit of work and begins to treat a multi-step workflow as the unit of work: intake, structuring, critique, revision, verification routing, and final drafting—each with explicit checkpoints and human approvals. In Level 3, the model is no longer just producing artifacts; it is moving artifacts through a governed pipeline.

This shift matters because it addresses the practical limit of Level 2: Level 2 can produce a verification plan, but it cannot enforce execution. Teams can generate beautiful lists of "things to verify" and then fail to verify them. Level 3 introduces mechanisms that make verification harder to skip. The workflow includes gates that block downstream steps until assumptions are registered and verification tasks are acknowledged, assigned, and, where feasible, completed or explicitly deferred with documented rationale.

Level 3 also introduces separation of duties. In a governed process, the same actor should not both generate a reasoning bundle and certify it as verified. Level 3 workflows therefore separate roles: a drafting pass, a skeptical review pass, a QA pass, and a human sign-off pass. This can be implemented with multiple model calls, but the critical point is not the number of calls; it is the checkpoints and the accountability attached to them. The model's outputs remain *Verification status: Not verified. Human review required.*until a human approves evidence and changes the status. In Level 3, the process itself becomes the artifact: each step is logged, each prompt is captured, each output is versioned, and each decision point is recorded.

Immutable logs become central at Level 3 because the organization now needs to prove not only what it thought, but how it behaved. The audit trail includes run manifests, prompt logs, risk logs, and artifact bundles for each matter or engagement. This is where "governance-first" becomes operational at scale. Level 3 is the beginning of agentic workflows, but the book's posture remains conservative: agents do not replace judgment; they enforce disciplined sequencing and recordkeeping. Human ownership is strengthened, not weakened.

The practical promise of Level 3 is therefore clear: the organization can execute the Level 2 discipline reliably, even under pressure, because the workflow is built to prevent shortcuts. The risks also increase: multi-step workflows create more complexity, more failure points, and more opportunities for silent drift. That is why Level 3 governance becomes stronger: checkpointing, separation of duties, QA routines, and immutable logs are not optional enhancements; they are the controls that make agentic workflows defensible.

---

**Artifact (Save This)**

**Level 2 exit criteria (ready to move to Level 3).**

1. The team reliably produces issue maps, alternatives, and tradeoffs without recommendations.
2. Assumptions are logged with owners and tests (assumption register is actually used).
3. Weakest-link assumptions are identified and drive verification priorities.
4. Reasoning bundles are saved with prompt logs, reviewer notes, and version history.
5. No decision is justified solely by AI-produced reasoning artifacts.

---

If a team can meet these exit criteria consistently, it has learned the core lesson of Level 2: structure is a tool for scrutiny, not a substitute for truth. With that discipline in place, the organization is ready for Level 3, where the same reasoning artifacts will be embedded into multi-step workflows that enforce checkpoints, preserve provenance, and make the process itself auditable.

# Bibliography

[1] National Institute of Standards and Technology (NIST). *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*. NIST AI 100-1, U.S. Department of Commerce, 2023.

[2] International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC). *ISO/IEC 23894:2023 — Information technology — Artificial intelligence — Guidance on risk management*. ISO/IEC, 2023.

[3] International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC). *ISO/IEC 42001:2023 — Information technology — Artificial intelligence — Management system*. ISO/IEC, 2023.

[4] International Organization for Standardization (ISO). *ISO 31000:2018 — Risk management — Guidelines*. ISO, 2018.

[5] International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC). *ISO/IEC 27001:2022 — Information security, cybersecurity and privacy protection — Information security management systems — Requirements*. ISO/IEC, 2022.

[6] European Parliament and Council of the European Union. *Regulation (EU) 2024/1689 of 13 June 2024 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act)*. Official Journal of the European Union, 2024.

[7] Organisation for Economic Co-operation and Development (OECD). *Recommendation of the Council on Artificial Intelligence (OECD/LEGAL/0449)*. Adopted 22 May 2019.

[8] Cabinet Office (United Kingdom). *Generative AI Framework for HM Government*. UK Government guidance, 2024.

# Chapter 3

# Agents

**Abstract.** Level 3 introduces *agentic* use of Generative AI in consulting and corporate strategy: not autonomy, but disciplined multi-step workflows that produce a sequence of intermediate artifacts under strict governance. The key change is that the unit of control is no longer a single memo or slide; it is the **workflow run** itself, with defined stages, mandatory human checkpoints, and a persistent audit trail. This chapter teaches how to design stage boundaries (inputs, outputs, acceptance criteria), enforce a facts-versus-assumptions schema at every step, and prevent downstream work from proceeding when required inputs are missing or when high-risk claims appear. Readers practice inserting review gates, defining stop conditions and escalation rules, and separating drafting, QA checking, and approval roles (even in small teams). The chapter's mini-cases (market entry, cost transformation, capital allocation, operating model redesign) show how the same strategy questions can be executed as governed workflows that generate inspectable artifacts: issue maps, hypothesis lists, risk registers, draft narratives, and verification question lists—each explicitly labeled *Not verified*. The result is faster production of draft bundles *without* outsourcing judgment, while improving defensibility through traceable process, version history, and documented human ownership.

---

**Artifact (Save This)**

**Scope disclaimer (required).** This chapter is an educational governance-first guide for management consulting and corporate strategy work. Outputs produced with AI are drafts only and must be reviewed by a qualified professional. The model may be persuasive and wrong. Do not treat outputs as factual, complete, or client-ready without verification and sign-off.

**Level alignment.** Level 3 introduces **agentic workflows**: multi-step processes that produce a sequence of intermediate artifacts with **mandatory human checkpoints** and **persistent logs**. No autonomous decision-making. No independent verification. Human judgment remains fully responsible.

## 3.1 Chapter overview: Level 3 in the maturity ladder

### 3.1.1 The five levels (why this ladder exists)

The maturity ladder in this book exists because Generative AI capability does not arrive all at once. It arrives in increments that *look* small from a distance ("the model is just a bit better"), but that create discontinuous jumps in organizational risk. The ladder is therefore not a marketing device. It is a governance device. Its purpose is to help consulting teams and strategy functions adopt useful capability *without drifting into uncontrolled delegation of judgment*, and to make explicit a principle that firms routinely forget under time pressure: as AI can do more, the organization must do more to constrain, document, review, and own what is produced.

Level 1 is the baseline: **Chatbots**. At this level, AI is a disciplined drafting and formatting assistant. It helps you turn rough notes into a coherent email, a memo outline, a meeting summary, or a slide storyline. The output can look impressively polished, but the governance posture must remain strict: it is a draft, it is not verified, and it is only as reliable as the facts you provided and the care with which you review it. The value proposition is speed and clarity, not truth or authority.

Level 2 introduces a distinct shift: **Reasoners**. Here, the model is not merely producing prose; it is being used to make thinking inspectable. You ask it to decompose issues, list alternatives, map tradeoffs, and stress-test assumptions. The output is no longer "a better paragraph," but a reasoning artifact that can be challenged. The central discipline at Level 2 is that reasoning must be *structured and explicitly provisional*. The model can help you see the shape of an argument, but it cannot certify that the argument is correct. Governance, therefore, focuses on separating facts from assumptions and forcing open questions into the foreground before a team drifts into persuasive logic.

Level 3 is the next jump: **Agents**. The term "agent" is overloaded and often abused. In this book it has a narrow, governance-first meaning: a multi-step workflow in which AI produces a *sequence* of intermediate artifacts, typically across multiple stages (intake, scoping, decomposition, drafting, QA checking, revision), with explicit human checkpoints that gate what happens next. Level 3 exists because real consulting and strategy work is inherently procedural. A one-page decision memo is not produced by a single clever paragraph; it is produced by a chain of decisions about scope, about what facts are known, about which assumptions are acceptable, about which stakeholders must be aligned, and about which risks must be acknowledged. Level 3 is where the organization admits that the work is a process and then governs that process.

Level 4 expands the focus: **Innovators**. At this stage the organization begins to build reusable internal assets: templates, playbooks, prompt libraries, workflow blueprints, evaluation checklists, and controlled variants for different contexts. The output is not only a deliverable, but also a reusable capability that can be deployed repeatedly. Because reuse scales impact, Level 4 requires stronger change control: versioning, test cases, release notes, and clear scope constraints. The organization must treat these assets as internal products with owners, not as casual documents.

Level 5 is the full organizational posture: **Organizations**. Here the maturity ladder becomes

an operating model. The firm has intake and classification, routing rules, access controls, audit trails, training, supervision, and an evaluation regime that is sustained over time. Level 5 is where governance stops being a set of individual best practices and becomes a system: a minimum standard for safe operation, with accountability, monitoring, and controlled evolution.

The ladder matters because it prevents a common failure mode in professional services: adopting the *most capable* tool with the *least capable* governance. Teams often try to jump from "chatbot drafting" to "autonomous agent" by accident, because they confuse output quality with process quality. This ladder is designed to stop that drift. It forces a team to earn capability by also earning the controls that make that capability defensible.

### 3.1.2 Why Level 3 exists

If Level 2 is about making reasoning visible, Level 3 is about making execution repeatable and auditable. This is not a cosmetic distinction. It is a recognition that consulting and corporate strategy work is rarely a single artifact. It is a chain of artifacts, produced across time, with dependencies. The chain typically begins with intake: what is the question, who is the decision-maker, what is the deadline, what is in scope and out of scope, what constraints exist (financial, regulatory, operational, political), and what facts are already known. It continues with scoping: what must be true for a recommendation to be viable, what uncertainties matter, what stakeholders must be aligned, what risks are unacceptable. It proceeds to decomposition: an issue map, a hypothesis set, a set of alternative pathways. Then drafting: the memo, the slide storyline, the workshop plan. Then QA: consistency checks, missing-field checks, contradiction scans, language that implies certainty. Then revision: human edits and integration. Finally, sign-off: a human owner takes responsibility for the deliverable and for the process by which it was produced.

This procedural reality creates an uncomfortable truth: *the biggest risks are often not in the final document, but in the process that produced it.* A team can produce a reasonable-looking memo that is fundamentally compromised because the scope was never agreed, because assumptions were silently introduced and then treated as facts, because a stakeholder constraint was ignored, or because a "QA pass" was mistaken for verification. Under time pressure, teams skip steps. Under excitement about AI, teams accelerate steps. In both cases the failure is not merely "bad text." The failure is "bad process at speed," where errors propagate quickly and become harder to detect because everything looks polished and consistent.

Level 3 exists to address this specific risk profile. It is the level at which the firm begins to treat AI usage as a workflow with gates, rather than as a tool used opportunistically. The practical reason is simple: once the model is asked to run multiple steps, it begins to feel like a junior consultant who can "just handle it." That feeling is precisely what must be governed. A multi-step workflow can create an illusion of due diligence: it can produce an issue tree, a risks section, a list of stakeholders, and a conclusion, all in one flow. But if the underlying facts were never confirmed, if the assumptions were never owned, and if the checkpoint was skipped, then the workflow is not diligence; it is automation bias.

In consulting and corporate strategy contexts, the danger is amplified by the nature of decisions. Many strategic decisions are irreversible or expensive to reverse: entering a market, exiting a product line, restructuring operations, allocating capital, changing an operating model. The cost of being "confidently wrong" is high. Moreover, the work is often communicated upward, where brevity and confidence are valued. A polished, decisive-sounding output can be more dangerous than a messy one, because it travels faster. Level 3 therefore insists on two things that may feel frictional but are essential: (1) the workflow must be designed to produce *interruptible* outputs, and (2) the workflow must be designed to produce *reviewable evidence* of how those outputs were generated.

The phrase "agentic workflow" in this book is therefore intentionally conservative. It does not mean the AI goes off and conducts research, makes calls, chooses a strategy, and returns with a plan. It means the AI executes a defined sequence of steps that the human team designed, under explicit constraints, and with mandatory gates. The model can draft an intake form, propose an issue map, generate alternatives, produce a first-pass memo, run a structural QA pass, and surface open questions. But at each stage it must label what is not verified and what must be checked, and at key points it must stop and hand off to a human owner.

The payoff is not magic. The payoff is discipline. A governed workflow reduces variability: different consultants and analysts produce artifacts in different ways; the workflow forces a common structure. A governed workflow improves traceability: if a partner asks, "Why did we say this?" the team can point to the stage where the claim was introduced, whether it came from a fact provided or an assumption introduced, and what review occurred. A governed workflow improves training: junior staff learn not only what a good memo looks like, but what a defensible process looks like. And, crucially, a governed workflow allows the organization to scale AI usage without scaling chaos. Without Level 3, scaling AI often means scaling untracked drafts, undocumented assumptions, and unreviewed outputs.

Level 3 also exists because it is the threshold at which organizations can meaningfully talk about auditability. At Level 1 and Level 2, you can save prompts and outputs, but the work is still largely artisanal: the sequence of steps is implicit. At Level 3, the sequence becomes explicit. Once the sequence is explicit, you can demand that it be followed. Once you can demand that it be followed, you can test whether it was followed. Once you can test, you can improve. This is the beginning of operational governance.

### 3.1.3 What changes at Level 3 (and what does not)

The most important change at Level 3 is that the "unit of work" shifts from a single artifact to a workflow run. That shift has several concrete implications.

First, work becomes stage-based. A Level 3 workflow is not "one prompt." It is a series of stages, each with a defined objective, required inputs, expected outputs, and acceptance criteria. This is not bureaucratic decoration. It is a control mechanism. When stages are defined, the workflow can enforce that certain information must be provided before a later artifact is produced. For example, a market entry decision memo should not be drafted before the scope is locked and the

decision-maker is identified, because the memo's structure and emphasis depend on those facts. A cost transformation workplan should not proceed past an initiative list if the baseline cost taxonomy is unknown, because the list will otherwise become generic and misleading. Stage discipline is how a workflow resists the model's tendency to "fill in the blanks."

Second, Level 3 introduces mandatory human checkpoints that gate progression. At Level 1 and Level 2, human review is required before anything leaves the team, but the internal flow can still be informal. At Level 3, informality becomes dangerous because the workflow can run fast. Checkpoints must therefore be explicit and enforceable. A checkpoint is not "someone should glance at this." A checkpoint is: a named human owner reviews specified fields (facts provided, assumptions, open questions, risk flags), makes a decision (approve, revise, reject), and that decision is recorded. The workflow must treat an unapproved checkpoint as a stop condition, not as a suggestion.

Third, Level 3 introduces explicit separation of duties as a design principle. In an ideal setting, the person (or role) drafting is not the person approving. The person performing a QA check is not the person who benefits from passing it. In small teams, the same person may necessarily wear multiple hats, but the hats must be explicit in the record. The point is not to create ceremony; the point is to prevent self-approval becoming invisible. Level 3 acknowledges a practical reality: the model can generate outputs that look correct, and humans can become complacent when they are tired. Separation of duties is a defense against that complacency.

Fourth, persistent logs and versioning become mandatory rather than optional. At Level 1 and Level 2, a team can get significant value simply by saving the final draft and the prompt that produced it. At Level 3, that is insufficient because intermediate artifacts matter. The organization must be able to reconstruct what happened across stages: what inputs were used, what prompts were applied, what outputs were produced, what changes were made, and who approved what. This is the minimum audit trail for professional defensibility. It is also the foundation for improvement: without logs, you cannot know whether a workflow is reliable or where it fails.

Fifth, stop conditions and escalation rules become first-class. A workflow that always proceeds is a workflow that will eventually proceed into error. Level 3 therefore requires explicit stop conditions, such as: missing required inputs, detection of high-risk claims (quantitative claims, external commitments, regulatory assertions), scope expansion beyond what was approved, or contradictions between facts and conclusions. When a stop condition triggers, the workflow must halt and escalate to a human owner. This is how the workflow prevents the model from "helpfully" completing the task by inventing what it does not know.

Sixth, Level 3 changes how teams think about QA. At Level 1, QA is largely human editing. At Level 2, QA may include structured reasoning checks. At Level 3, QA often includes an AI-performed *non-verifying* check stage: completeness checks, structure checks, contradiction scans, and language-risk scans (for example, flagging phrases that imply certainty or verification). This is useful, but it introduces a new risk: teams may confuse QA with verification. Level 3 must therefore be explicit that AI QA is not truth validation. It is only a consistency and structure tool. Verification remains a human responsibility, often requiring external evidence.

These changes are substantial, but they are paired with what does *not* change, and the non-changes are as important as the changes.

AI still does not verify. The workflow may include a stage called "verification planning," but that stage produces a list of questions and evidence needs; it does not produce verified facts. If the workflow includes research steps, those steps must be tightly governed and usually rely on human-provided sources or explicitly approved data. The model must not fabricate citations, benchmarks, market statistics, or "industry norms." Any quantitative claim must either be provided by the human or be flagged as an open question to verify.

AI still does not recommend. A workflow can enumerate options and map tradeoffs, but the decision remains a human act. In many consulting contexts, "recommendation" is the core value. Level 3 deliberately resists the temptation to let AI "choose" because choosing is not merely reasoning; it is judgment under accountability. The model can help structure the argument for an option, but it cannot own the option.

AI still is not autonomous. The workflow may feel agentic, but it is bounded by design. It has defined stages and gates, and it is allowed to stop. A Level 3 workflow is therefore *anti-autonomy* in spirit: it is a system that makes autonomy difficult by requiring approvals and by preserving human ownership.

Finally, accountability is not outsourced. This is not a philosophical statement; it is a practical one. In consulting and corporate strategy, accountability exists in the form of partner sign-off, executive sign-off, board accountability, and sometimes regulatory scrutiny. Level 3 strengthens the evidence trail that supports accountability, but it does not shift accountability away from humans. If anything, Level 3 makes accountability more explicit, because it forces the organization to name owners and to record decisions.

The net effect is that Level 3 is where AI stops being "a tool you use" and becomes "a process you run." That process can be safe and valuable, but only if it is designed with governance-first principles: facts separated from assumptions, checkpoints that cannot be skipped, logs that cannot be lost, and stop conditions that prevent the model from manufacturing certainty. If Level 1 taught teams to draft with discipline, and Level 2 taught teams to reason with discipline, Level 3 teaches teams to *execute* with discipline. The organization learns to scale work without scaling uncontrolled delegation.

---

**Artifact (Save This)**

**Mental shift at Level 3.** At Level 1 and Level 2, the unit of work is an artifact (email, memo, issue tree). At Level 3, the unit of work is a **workflow** that reliably produces a **bundle of artifacts** with checkpoints, traceability, and review evidence.

---

## 3.2 Mental model: what an "Agent" is in consulting work

### 3.2.1 Useful abstraction

In consulting and corporate strategy, the word "agent" is both powerful and dangerous because it invites a familiar analogy: a junior team member who can take a workstream, run with it, and come back with a coherent output. That analogy is precisely the one that creates governance failures. A safer abstraction is intentionally narrower: at Level 3, an "agent" is a **workflow operator** under constraints. It is a mechanism for executing a predefined sequence of steps and producing intermediate artifacts in consistent formats, so that humans can review, accept, revise, or reject those artifacts at explicit checkpoints. In other words, it is not a substitute for professional judgment; it is a substitute for process inconsistency.

This distinction matters because consulting work is rarely just thinking or just writing. It is thinking *as a process*, where the final deliverable is the last link in a chain of intermediate choices. Before the executive memo exists, the team has already made a set of decisions about what the question is, what the scope is, what constraints matter, what alternatives exist, what data is missing, and what risks must be highlighted. In a strong team, those intermediate choices are made explicitly; in a rushed team, they are made implicitly, often by whoever speaks last in the room. Level 3 tries to institutionalize the strong-team behavior: it turns intermediate choices into intermediate artifacts, each produced with a stable schema and each routed through a checkpoint.

A Level 3 "agent," therefore, is best understood as a disciplined procedure that creates a record. It takes an intake (objective, audience, constraints, sensitivity), produces a scoping statement, produces a decomposition artifact (issue map and hypotheses), produces a drafting artifact (memo or slide storyline), produces a structural QA artifact (consistency and completeness checks), and produces a handoff artifact (what is verified, what is not, what must be verified, who owns which decisions). The "agent" may be implemented as multiple prompts or multiple internal roles, but the defining property is not the number of prompts; it is the existence of **stage boundaries and gates**. In a Level 3 design, the workflow is allowed to stop, and it must stop when inputs are missing, when high-risk claims appear, or when the scope is being exceeded.

This is why "workflow operator" is the right abstraction. Consulting teams already operate within workflows, even when they pretend they do not. An engagement has an intake, a diagnostic, a hypothesis phase, a workplan, synthesis, and communication. A strategy function has a calendar cycle, a set of recurring decision memos, an approval process, and an escalation path. Level 3 does not invent process; it makes process explicit, and then governs it. When process is explicit, it can be inspected and improved. When process is implicit, it can only be hoped for.

The "operator under constraints" framing also clarifies what the system is *for*. It is not for brilliance. It is for repeatability. Most strategic failures inside organizations are not caused by the absence of brilliant ideas; they are caused by inconsistent execution of basic discipline: failure to clarify scope, failure to surface assumptions, failure to align stakeholders, failure to document tradeoffs, failure to notice that a claim rests on a missing data point. A Level 3 agent is a mechanism

to reduce those failures by making the workflow insist on them. The value is that the workflow does not get tired. It does not "forget" to include an assumptions register. It does not "forget" to list open verification questions. It does not quietly move forward when the decision-maker has not been identified. It can be designed to be stubborn in precisely the ways a well-governed consulting process should be stubborn.

A practical way to internalize this abstraction is to treat the agent as a "structured conveyor with inspection stations." The conveyor moves work from stage to stage. The inspection stations are the human checkpoints. The conveyor can produce intermediate artifacts quickly, but the stations decide whether those artifacts are acceptable inputs to the next stage. In this mental model, speed is not the objective; controlled flow is the objective. A workflow that moves quickly into error is not an improvement; it is a multiplier of risk. Therefore, the design goal is not maximal throughput; it is maximal defensibility per unit of throughput.

The workflow operator abstraction also forces clarity about what the system "knows." At Level 3, the agent is not discovering truth; it is transforming and organizing inputs. It can restructure your notes into an executive-ready narrative, but it cannot confirm that the facts are correct. It can list plausible alternatives, but it cannot guarantee that the list is complete. It can run a consistency scan, but it cannot validate the external world. This is not a limitation to apologize for; it is the foundation of safe use. Once a team accepts that the agent is an operator, not an oracle, it becomes natural to treat outputs as drafts and to treat verification as a separate human-owned workstream.

This operator abstraction also provides a clean way to define internal roles within the workflow. A common pattern is to split the work into at least three conceptual roles, even if the same model is used for all three: (1) **Drafting role**, which produces stage outputs; (2) **QA role**, which checks structure, completeness, contradictions, and language that implies certainty; and (3) **Gatekeeper role**, which is the human who approves progression. The agent can perform the first two roles, but the third must remain human. The agent can suggest a stage plan, but a human must approve it. The agent can propose an issue map, but a human must confirm that it matches the actual decision context. The agent can draft a memo, but a human must decide whether the memo is appropriate to send, and must own the risks of sending it.

Finally, the workflow operator abstraction helps teams avoid a subtle but important confusion: the difference between *performing steps* and *owning outcomes*. In consulting, outcomes are accountable: recommendations lead to investments, reorganizations, and public commitments. Level 3 allows AI to perform steps, but it does not allow AI to own outcomes. The workflow is therefore designed so that the output that matters most (the decision) is never produced as a model-owned statement. Instead, the workflow produces a bundle that supports a human decision: structured alternatives, tradeoffs, assumptions, and a verification plan. The human decides, and the record shows how the decision was supported.

### 3.2.2   Dangerous misconception

The most common misconception about agents is the simplest one: **agent = autonomy**. In popular discourse, "agent" is often used to describe a system that can independently set goals, acquire information, take actions, and iterate until an outcome is achieved. That is not what Level 3 means in this book, and it is not what a governance-first consulting practice can safely permit in most contexts. Autonomy is not a feature you sprinkle onto professional judgment; it is a liability that must be justified, bounded, and supervised. In strategy work, where the cost of confident error is high, autonomy is rarely the correct default.

The second misconception is equally seductive: **multi-step = correctness**. A workflow with multiple steps can feel like a sign of rigor. If the system produces an intake summary, then an issue map, then a tradeoff table, then a risk register, and then a conclusion, it looks as though due diligence has occurred. But the presence of steps is not evidence that the steps were grounded in verified facts. A multi-step process can be an elegant machine for manufacturing consistency around a mistaken premise. If the workflow begins with an incorrect "fact provided" or an unowned assumption, every downstream artifact will inherit it. The workflow then becomes a coherence amplifier: each stage reinforces the same premise, producing a bundle that looks internally aligned and therefore "credible," even though it is aligned around an error.

This error propagation dynamic is one of the defining risks at Level 3. At Level 1, hallucinations and invented facts are dangerous, but they are often localized to a paragraph. At Level 2, persuasive logic and hidden assumptions are dangerous, but they can be challenged within a reasoning artifact. At Level 3, the risk is that a single error becomes institutionalized across a bundle and then travels through the organization as a packaged narrative. The more artifacts the workflow produces, the more surface area exists for the error to appear "supported." The team may find it harder to identify the root cause because the error is no longer a single sentence; it is woven into the structure of multiple outputs.

A related misconception is that "QA" can replace verification. Many Level 3 workflows introduce an AI QA stage that checks for completeness, contradictions, and risky language. This is useful, but it creates the temptation to treat QA as validation. The model can detect that a memo contradicts itself, but it cannot detect that the memo contradicts reality unless the contradiction is explicitly represented in the inputs. A workflow that contains a QA stage is therefore not necessarily safer; it is safer only if the organization remains disciplined about what QA is. QA is a structural check. Verification is an evidentiary act performed by humans.

Another misconception is that workflows reduce responsibility because they reduce discretion. Teams may assume that if the workflow is standardized, then the workflow "owns" the output. This is how governance slowly fails: responsibility becomes diffused across process. The correct posture is the opposite. Standardization makes responsibility *more* visible, not less. If a workflow is standardized, then the organization can precisely ask: who approved scope, who approved assumptions, who approved the final output, and what evidence was used? The workflow makes those questions easier, which increases accountability. Therefore, a Level 3 system must be designed

so that human ownership is explicit at checkpoints, not implicit in hindsight.

There is also a misconception specific to professional services culture: that agents are a substitute for training. Leaders may hope that a workflow can compensate for junior staff who do not yet know how to write an issue map or structure a memo. In reality, a workflow can produce a plausible artifact, but it cannot teach judgment unless humans are forced to review and explain. If juniors simply accept outputs, the workflow becomes a shortcut around learning. A governance-first Level 3 implementation therefore treats the workflow as a training scaffold: juniors learn by seeing the structure, but they must still explain the assumptions, defend the tradeoffs, and answer the open questions. The workflow should increase professional discipline, not replace it.

In consulting work, the most dangerous misconception is that agentic workflows create "analysis." A model can generate reasoning-looking content with impressive fluency. It can produce a market entry narrative with crisp bullet points and a confident tone. Teams can confuse that fluency with analysis because in consulting, analysis often ends up being summarized as bullet points. The distinction is that real analysis has evidence, traceability, and a chain of reasoning tied to verified inputs. A Level 3 workflow can help create the structure for analysis, but it cannot create the evidence. Therefore, any workflow that allows the model to invent numbers, benchmarks, or external claims is not a Level 3 workflow; it is a risk engine.

Weak stage boundaries worsen every misconception. If stages are vaguely defined ("do the market entry analysis"), the model will fill in blanks, and the team will treat the output as "what the workflow produced." If checkpoints are weak ("looks fine"), the team will proceed without confronting open questions. If stop conditions are absent, the workflow will proceed into uncertainty and manufacture specificity. If logs are not saved, the team will be unable to reconstruct how a claim entered the bundle, and the organization will default to trusting the most polished artifact. Each of these weaknesses turns the workflow into a credibility generator rather than a governance mechanism.

To counter these misconceptions, Level 3 requires a deliberate linguistic discipline. Teams must avoid describing the agent as "doing the work" in an ownership sense. The agent is "producing a draft stage output" or "generating a candidate structure." Teams must avoid describing QA as "checking the facts." QA is "checking internal consistency and completeness." Teams must avoid describing the workflow as "giving the answer." The workflow is "producing a bundle to support a human decision." These are not semantic games; they are governance cues that prevent organizational drift.

The final misconception worth naming is the idea that a workflow is safe because it is internal. Many strategy artifacts remain inside a firm, and teams may feel that internal drafts can be looser. In practice, internal drafts are often more dangerous than external ones because they shape decisions before verification occurs. A flawed internal memo can lead to a flawed executive belief, which then drives resource allocation. Once executives internalize a narrative, it becomes hard to dislodge. Therefore, internal use does not reduce the need for governance. It increases the need for disciplined labeling (*Verification status: Not verified. Human review required.*), explicit open questions, and

checkpointed review, because internal drafts have direct influence on decisions.

### 3.2.3 Definition of a good Level 3 output

A good Level 3 output is not a single document. It is a **governed workflow run** that produces a bundle of artifacts, along with the minimum evidence required to defend how those artifacts were produced. The bundle is designed to make work reviewable, interruptible, and improvable. It is intentionally incomplete where facts are missing, and it is explicitly provisional where assumptions are introduced. It does not attempt to hide uncertainty; it organizes uncertainty so that humans can resolve it.

The first property of a good Level 3 output is that it is **stage-structured**. Each stage produces an intermediate artifact with a clear objective, and each artifact is saved. A typical Level 3 workflow might include: Stage 0 (intake and classification), Stage 1 (scope lock and constraints), Stage 2 (issue map and hypotheses), Stage 3 (alternatives and tradeoffs), Stage 4 (draft deliverable), Stage 5 (QA pass), Stage 6 (red-team critique), Stage 7 (human synthesis and sign-off). The exact stages vary, but the principle does not: each stage has defined inputs and outputs, and the workflow does not treat the final deliverable as the only thing worth saving.

The second property is that a good Level 3 output maintains **explicit unknowns**. Every stage output must separate facts provided, assumptions introduced, and open questions. This is not merely a formatting requirement. It is the central defense against persuasive coherence. A workflow that produces a beautiful memo without open questions is not being helpful; it is being dangerous. A strong Level 3 output makes missing information visible early, and it routes the resolution of unknowns to a human owner. It is acceptable for an output to be incomplete; it is unacceptable for an output to be complete-looking when it is unverified.

The third property is that a good Level 3 output includes a consolidated **assumptions register**. Assumptions introduced at any stage must be gathered into a single register that is easy to review. The register should indicate (at minimum) the assumption statement, the owner (human), the stage in which it was introduced, and whether the assumption is approved, pending, or rejected. In many workflows, assumptions should also be tagged for risk. For example, assumptions that influence quantitative claims, external commitments, regulatory positions, or stakeholder responses should be tagged as high-risk. The reason is practical: not all assumptions are equally dangerous. A harmless assumption about slide formatting does not deserve the same gate as a critical assumption about market growth. The register allows the team to focus attention where it matters.

The fourth property is that a good Level 3 output includes a **verification plan**. In governance-first practice, verification is not an afterthought; it is a workstream. The workflow should produce a list of verification questions and the evidence needed to resolve them. This plan does not require the model to perform the verification. It requires the model to help structure what must be verified, by whom, and in what order. In a market entry context, verification questions might include: what is the true size of the addressable market under the agreed scope, what regulatory approvals are required and under what timeline, what competitor responses are historically observed in similar

entries, what internal capabilities are prerequisites, and what financial constraints apply. The verification plan turns "we should check this" into a trackable list, which helps prevent the team from accidentally sending a memo that relies on unchecked claims.

The fifth property is that a good Level 3 output includes **human review notes and a sign-off path**. This is where Level 3 becomes meaningfully different from Level 2. A Level 2 reasoning artifact can be excellent and still fail in practice if there is no record of who accepted the assumptions and who decided to proceed. Level 3 requires that checkpoints are recorded. The workflow run should indicate: who approved scope, who approved assumptions, who reviewed the draft deliverable, what changes were made, what remains open, and who owns the final decision to communicate the output. Even if the process is lightweight, the ownership must be explicit, because that is what makes the workflow defensible.

The sixth property is that a good Level 3 output is **safe by design under deadline pressure**. This may be the most important criterion. In consulting, the workflow will often be run when people are tired and the deadline is close. A good Level 3 design anticipates this and refuses to progress silently. If required inputs are missing, the stage output must say so and must trigger a stop condition. If the draft includes a quantitative claim not present in facts provided, it must be flagged as an open question rather than presented as a confident statement. If the scope is ambiguous, the workflow must force a scope lock before proceeding. The test of a Level 3 workflow is not whether it works in a calm environment; it is whether it prevents reckless speed in a pressured environment.

A good Level 3 output also respects a conservative "minimum necessary" stance on content. The workflow should encourage teams to provide only the facts required to structure the work, and to redact sensitive details. This is not merely a confidentiality posture; it is a quality posture. When teams paste large quantities of raw sensitive content, they often reduce their own ability to review. A governance-first workflow encourages small, explicit facts and explicit unknowns.

Finally, a good Level 3 output is designed for improvement. Because the workflow is logged, the organization can learn from failure. If the workflow repeatedly produces open questions in the same area, the team can adjust the intake stage to capture that information earlier. If the workflow repeatedly produces risky language that implies certainty, the QA stage can be strengthened to flag those phrases. If the workflow repeatedly generates overly generic alternatives, the scoping stage can be revised to tighten constraints. This is one of the practical benefits of Level 3: it creates the foundation for Level 4 innovation, where workflows become reusable assets with controlled evolution. But even at Level 3, the improvement loop is already present: a workflow run is not only a deliverable; it is data about the reliability of your process.

In short, a good Level 3 output is a bundle that allows a professional to answer three questions with integrity: (1) What do we know, and how do we know it? (2) What are we assuming, and who owns those assumptions? (3) What must be verified before we treat this as decision-ready? If the bundle supports those questions, it is a governed workflow run. If it does not, it may be a polished set of drafts, but it is not Level 3 in the governance-first sense.

---

**Artifact (Save This)**

**Non-negotiable rule. Facts are not assumptions.** Facts must be provided or verified; assumptions must be stated, owned, and testable.

**Required fields for every Level 3 stage output (minimum).**

a. facts_provided (approved inputs only)
b. assumptions (explicit, owned; may block downstream stages)
c. open_questions (what must be verified by humans)
d. stage_draft_output (the artifact produced in this stage)
e. verification_status: `"Not verified"`

---

## 3.3   What Level 3 can do, and what it cannot do

The point of the CAN / CAN'T section is not to restate common sense. In professional environments, common sense is what disappears first under deadline pressure, stakeholder anxiety, and the seductive smoothness of AI outputs. Level 3 introduces workflows that can run quickly and produce many artifacts in sequence, which means the risk profile changes: the organization is no longer merely managing the risk of a single flawed paragraph, but the risk of *a flawed process that produces a coherent bundle.* This section therefore defines Level 3 capabilities in a way that is operational: what a governed workflow may do, what it must do, and what it must refuse to do.

A helpful way to think about Level 3 is that it formalizes a specific division of labor. AI performs transformations and consistency checks on structured inputs; humans perform judgment, ownership, verification, and approval. If a workflow blurs that division, it ceases to be governance-first. In practice, many teams drift into blur because "the workflow already did the steps." The purpose of the CAN / CAN'T list is to prevent that drift by anchoring what the workflow is permitted to output, what it is permitted to imply, and what it must explicitly flag as not verified.

### 3.3.1   CAN: workflow execution under governance

A Level 3 system *can* run a predefined multi-step workflow that produces multiple intermediate artifacts, provided that the workflow is explicitly scoped, checkpointed, and logged. That sentence contains the core constraints. "Predefined" means the organization has decided in advance what stages exist, what each stage does, and what inputs are required before the stage can run. "Multi-step" means the system's output is not only the final deliverable but also the intermediate artifacts produced along the way. "Under governance" means the workflow does not run as an uninterrupted stream; it runs with gates, records, and stop conditions.

This is often the first capability that teams experience as qualitatively different from Level 2. At Level 2, the model may produce an issue map or a tradeoff table, but the human orchestrates the overall sequence informally. At Level 3, the sequence itself becomes a design object. The workflow is structured so that it reliably produces, for example: an intake summary, a scope lock

statement, an issue map, a hypothesis set, an alternatives table, a risk register, a draft memo, a QA report, and a handoff note. Each artifact is saved, each artifact is labeled *Verification status: Not verified. Human review required.*, and each artifact is associated with the stage that produced it. The value is repeatability: different team members can run the same workflow and produce comparable structures, which reduces variation and reduces the "personal style risk" that often governs how consulting work is executed.

However, a workflow can only be "run" safely when it is allowed to stop. In governance-first Level 3, a workflow is designed to halt when required inputs are missing or when risk flags are triggered. For example, if the intake stage is missing the decision-maker, the workflow may not proceed to drafting a decision memo, because the memo's form and tone depend on who is deciding. If the workflow detects that the draft includes quantitative claims not present in facts provided, it must halt or at minimum flag them as open questions requiring human verification. The system can run quickly, but it cannot run blindly.

A Level 3 system *can* apply structured templates consistently across stages. This is where the workflow operator abstraction becomes practical. Consulting work benefits from standardized structures: issue maps, hypotheses, alternatives lists, tradeoff matrices, risk registers, stakeholder maps, and decision memos. When the model is used as a workflow operator, it can produce these structures with consistent headings and consistent fields. Consistency is not merely aesthetic. It is a control. When every stage output has the same schema (facts provided, assumptions, open questions, stage draft output, verification status), reviewers know where to look. They can quickly compare two runs. They can quickly detect missing sections. They can more easily identify where assumptions were introduced. Template consistency is therefore an accelerant for review discipline.

This capability is especially valuable in consulting environments where multiple workstreams feed a synthesis. If each workstream uses an idiosyncratic structure, the synthesis lead spends time normalizing before thinking. A Level 3 workflow can normalize by design. For example, every workstream can produce an issue map that uses the same decomposition conventions, every risk register can include the same categories, and every draft narrative can include a "what would change my mind" section. The model is effective at this kind of structured repetition. Governance-first practice harnesses that strength to reduce the cognitive load on humans, so humans can spend their limited time on judgment rather than formatting.

A Level 3 system *can* perform **non-verifying** QA passes. This is a subtle but important capability. The model can be used as a consistency checker that looks for internal contradictions, missing required fields, ambiguous pronouns, scope drift, and language that overstates certainty. For example, it can detect that a draft memo claims "we have confirmed" when no confirmation step exists in the workflow run. It can detect that the draft introduces a market size number that is not present in facts provided. It can detect that the conclusion implies a recommendation when the workflow is meant to be non-recommending. It can detect that the risk register is missing execution risks or missing second-order effects. These checks can meaningfully improve quality and reduce superficial errors.

But the governance emphasis must be explicit: QA is not verification. A non-verifying QA pass does not check whether the market size is correct; it checks whether the market size is *grounded in provided facts or flagged as an open question.* It does not check whether a regulatory claim is true; it checks whether the claim is accompanied by a verification question and a cautious posture. QA at Level 3 is best described as "defensive editing with structure," not as fact-checking. If an organization treats QA as verification, Level 3 becomes dangerous because it provides the illusion of validation without the reality of evidence.

A Level 3 system *can* produce handoff notes for human reviewers. This capability is often overlooked, but it is one of the most practically valuable features of an agentic workflow. In a professional setting, the limiting resource is not the model's ability to produce drafts; it is human review bandwidth. Handoff notes are designed to make review tractable by telling the human reviewer what matters: what changed between iterations, what assumptions were introduced, what open questions are blocking, what risk flags are present, and what the workflow recommends as next verification steps. These notes turn review from an unstructured reading task into a targeted inspection task.

Good handoff notes also preserve accountability. They explicitly state what the model did and did not do, and they remind the reviewer that the output is not verified. They can include suggested "review checks," such as: confirm that facts provided are complete, confirm that assumptions are acceptable, confirm that the output does not imply verification, confirm that stakeholder constraints are reflected. In other words, handoff notes help humans behave like professionals even when they are busy. That is the purpose of a governed workflow: not to remove humans, but to make humans more reliable.

There is one more "CAN" that is implied by the list but worth stating explicitly: at Level 3, AI can help the organization *enforce discipline* by making certain behaviors difficult. For example, if the workflow requires that every stage output include a verification status field, then the model will consistently include it. If the workflow requires that assumptions be explicitly listed, then assumptions will not remain hidden as easily. If the workflow includes a stop condition for missing inputs, then the model will refuse to proceed. This is a practical form of governance. It is not a policy document; it is a behavior embedded in the workflow.

### 3.3.2 CAN'T: things you must not delegate to Level 3

The most important "CAN'T" is the simplest: a Level 3 workflow must not make the decision, select the recommendation, or present a conclusion as "the answer." This is not because the model cannot produce a recommendation. It can, and it often will, because it is trained to be helpful. The reason is governance: recommending is a form of professional judgment that carries accountability. In consulting and corporate strategy, recommendations are not merely logical outputs; they are commitments that shape investments, reputations, careers, and sometimes legal obligations. A workflow can support a recommendation by producing alternatives and tradeoffs, but the act of choosing must remain human-owned, documented, and accountable.

This boundary is especially important because Level 3 workflows create "completion pressure." When a workflow has produced a bundle—issue map, tradeoffs, risks, narrative—people may feel that it is "ready" and therefore that the recommendation is the natural next step. But the model does not know what the organization can actually execute, what political constraints exist, what risk appetite is acceptable, or what values should govern the decision. Those elements are not mere facts; they are contextual judgments. Even if the workflow includes them as "facts provided," the final choice is still a human act. Therefore, a governance-first Level 3 design often uses language that avoids recommendation framing. It uses "options," "considerations," "tradeoffs," and "questions to resolve," and it requires the human owner to write the final decision statement.

The second "CAN'T" is that the workflow must not verify facts, benchmarks, citations, or external claims without human validation. This is the boundary that protects the organization from hallucinated authority. In consulting work, there is a constant temptation to include numbers: market size, CAGR, competitor share, cost benchmarks, headcount ratios, implementation timelines. These are precisely the types of claims that models can produce with fluent confidence and that busy executives may accept without checking. A Level 3 workflow must therefore be explicit: unless the number is provided as a fact, it must be treated as an open question to verify. The workflow can propose a place-holder ("TBD") and can propose what evidence to seek, but it cannot invent the evidence.

This also applies to "citations." In a governance-first posture, a model should not fabricate references or "industry reports" to sound credible. If the workflow needs references, the references must be provided or must be retrieved through controlled mechanisms that are explicitly outside the model's generative step (and, in many organizations, are human-performed). The model can help draft a verification plan: "find authoritative sources for X," "confirm with internal finance," "check regulatory guidance," but it cannot label claims as verified.

The third "CAN'T" is that the workflow must not proceed past checkpoints without explicit approval. This is the control that distinguishes Level 3 from informal chaining. The temptation to skip checkpoints is strong because checkpoints are friction, and friction feels expensive when a deadline looms. But the absence of checkpoints is what allows the workflow to become autonomous by default. A Level 3 workflow is designed to prevent "auto-approval." If the design allows stages to proceed without a recorded approval, the workflow is not governed; it is merely automated.

In practice, this means the workflow must define what approvals are required and what constitutes approval. Approval should be a deliberate act by a named human owner, recorded in the workflow run bundle. It can be lightweight (e.g., initials and date), but it must exist. The organization should also define which checkpoints are "hard gates" (must stop) and which are "soft gates" (can proceed with warnings). Governance-first posture recommends more hard gates than teams initially prefer, because the cost of one unapproved assumption propagating through the workflow can exceed the cost of delay. A workflow that runs fast but cannot prove it was reviewed is a workflow that will eventually produce defensible-looking errors.

The fourth "CAN'T" is that the workflow must not expand scope beyond the intake and scoping

gate. Scope creep is one of the most common consulting failure modes even without AI. With AI, it becomes easier because the model will happily answer adjacent questions, and it will often do so in a way that appears helpful. In strategic work, adjacent questions can be high-risk: a scope-limited request for an internal outline can drift into client-ready commitments; a request for an alternatives list can drift into a recommendation; a request for qualitative considerations can drift into invented quantitative claims. Therefore, Level 3 workflows require explicit scope lock early, and they require the model to treat anything outside that scope as an open question or an out-of-scope item.

Scope control must also be enforced at the language level. A workflow may be scoped to produce a draft decision memo, but the memo must not claim that analysis has been performed beyond what was provided. It must not imply that stakeholder alignment has occurred. It must not imply that data has been collected. If the workflow is scoped to produce a workshop agenda, it must not present workshop outcomes as known. In other words, scope control is not merely a project management concept; it is a content control. The output must match the scope of the workflow run.

Beyond these four "CAN'T" items, there are additional non-delegable boundaries that are implied by the governance-first posture and that Level 3 makes more urgent.

A Level 3 workflow must not allow the model to "decide by phrasing." Even if the workflow does not explicitly say "recommend," a model can recommend implicitly through tone and ordering: listing one option as obviously superior, using language like "clearly," "must," "should," or framing other options as weak. This is a subtle form of de facto decision-making. A governed workflow must therefore include language discipline, and the QA stage should flag recommendation-like phrasing. The human reviewer should be trained to detect when the model is smuggling judgment through rhetorical confidence.

A Level 3 workflow must not allow the model to collapse uncertainty. Models tend to produce completeness. If a workflow asks for an issue map and the input facts are sparse, the model will still produce a full map, which can mislead reviewers into thinking the map is grounded. Governance-first practice treats this as a risk: completeness can mask missing evidence. Therefore, the workflow must force explicit unknowns and must treat missing information as an output, not as a failure. In other words, the workflow must be designed so that "I need more facts" is a valid stage output. If the workflow punishes missing facts by forcing the model to proceed anyway, it will create invented specificity.

A Level 3 workflow must not allow "QA theater" to substitute for skeptical review. A strong QA stage can become a crutch: teams may rely on the model to catch issues and then reduce their own scrutiny. But the model's QA is limited to what is in the text. It will not notice that a critical stakeholder is missing if the stakeholder was never mentioned. It will not notice that an internal constraint is violated if the constraint is not present. Therefore, the organization must treat AI QA as a helpful assistant, not as a gatekeeper. The gatekeeper remains human.

A Level 3 workflow must not degrade confidentiality and data governance. Because Level 3 produces more artifacts, it increases the temptation to paste more raw content "so the workflow can handle it." This is dangerous both from a confidentiality perspective and from a reviewability

perspective. Governance-first Level 3 requires minimum-necessary input and strong redaction discipline, as well as clear policies on what can be pasted into the model. The workflow should encourage placeholders and local mappings, not raw client identifiers. This is not merely compliance; it is a control that prevents accidental leakage into logs and artifacts.

Finally, a Level 3 workflow must not become a substitute for professional accountability. The workflow can produce an audit trail, but the audit trail is not a shield. In fact, it is often a spotlight. If the workflow shows that a checkpoint was skipped, that record will harm, not help. Governance-first Level 3 therefore requires that the workflow is designed to be usable in reality: checkpoints must be simple enough to execute, logs must be easy to save, and roles must be clear. Governance that cannot be followed under real conditions is governance that will be silently ignored.

Putting it all together, Level 3 "CAN" is about transforming the consulting process into a repeatable, inspectable workflow that produces consistent intermediate artifacts and helps humans review more effectively. Level 3 "CAN'T" is about preventing the workflow from slipping into autonomy, verification theater, or rhetorical decision-making. The boundary is not a theoretical line. It is an operational line enforced through stage schemas, gates, stop conditions, logs, and language discipline.

In the next sections, the chapter will move from these boundaries to design patterns: how to define stages, how to place checkpoints, how to construct stop conditions, how to separate drafting from QA from approval, and how to save the workflow run bundle. The goal is to give teams a practical playbook: not "use agents," but "run governed workflows" in a way that makes outputs faster to produce, safer to use, and easier to defend.

## 3.4 The Level 3 governance wrapper (minimum controls)

Level 3 is where governance stops being a set of good intentions and becomes an engineered constraint system. At Level 1, a team can remain reasonably safe if it remembers to label outputs as drafts, separates facts from assumptions, and performs basic human review. At Level 2, safety improves when reasoning is made explicit and structured, because hidden assumptions and persuasive logic become easier to see. Level 3 changes the game again because the organization is no longer governing a single artifact. It is governing a *workflow* that can produce many artifacts quickly, in a sequence, with dependencies. Speed is not the enemy; ungoverned speed is. The Level 3 governance wrapper exists to ensure that the workflow cannot silently transform AI capability into de facto autonomy, cannot silently transform internal consistency into "verification," and cannot silently transform incomplete inputs into complete-looking outputs.

The most important principle behind the Level 3 wrapper is that every additional step in a workflow is an additional opportunity for error propagation. A single invented number or a single unowned assumption can appear in Stage 2 and then be repeated, rephrased, and embedded into Stage 4 and Stage 6 until it feels "established." This is why Level 3 governance is not primarily about better prompts; it is about **workflow design**: stage boundaries, gates, schemas, logs, stop

conditions, and role separation. The wrapper is a minimum control set that turns a multi-step process into something you can defend under scrutiny.

There are three reasons the wrapper must be minimal yet strict. First, consulting teams operate under time pressure. A control framework that requires elaborate ceremony will be ignored. Second, at Level 3 the volume of output increases: multiple stage outputs per run, potentially across multiple workstreams. Controls must be scalable and repeatable. Third, the organization is training itself. Early Level 3 adoption is a discipline-building phase. Controls must be simple enough to execute, but strict enough to actually shape behavior.

The wrapper can be summarized as: **define the workflow, gate the workflow, log the workflow, and label the workflow.** Each component addresses a predictable failure mode.

**Define the workflow.** A workflow that is not defined will be invented by the model. When the model invents a workflow, it invents not only steps but also implied evidence ("we reviewed," "we confirmed," "the analysis shows"). Therefore, Level 3 requires explicit stage definitions with acceptance criteria. Each stage must specify: what it takes as inputs, what it is allowed to produce, and what constitutes a valid output. This prevents the model from "helpfully" performing tasks that the organization did not authorize (like recommending, verifying, or expanding scope).

**Gate the workflow.** A defined workflow is still dangerous if it can proceed without approval. Gates prevent checkpoint skipping, which is one of the most common real-world breakdowns. A gate is not a suggestion to review; it is a hard stop. A workflow cannot proceed to the next stage unless the gate is explicitly cleared by a named human owner. This is how the organization prevents "auto-approval" behavior and preserves accountability.

**Log the workflow.** Without logs, governance becomes theater. When a partner or executive asks "How did we get here?", the team must be able to reconstruct what was provided, what was assumed, what was generated, what was changed, and who approved. Logs and version history are therefore not optional documentation; they are the core of defensibility. They also enable improvement. You cannot fix a workflow you cannot observe.

**Label the workflow.** Labeling is not cosmetic. In professional settings, the label "Not verified" is a control because it prevents downstream users from treating the output as fact. At Level 3, labeling must be embedded into every stage output because a workflow produces many outputs and any one of them may travel. A risk register can be copied into an email; an issue map can be pasted into slides; a draft narrative can become the basis for a leadership meeting. Every artifact must therefore carry the same warning posture.

These principles are implemented through specific minimum controls: stage definitions and acceptance criteria, human gates, schema enforcement, prompt logs, version history, stop conditions, role separation, and explicit "Not verified" labeling. The following risk-and-control framing captures why these controls are mandatory, not optional.

---

**Risk & Control Notes**

**Capability.** Checkpointed workflow execution for consulting artifacts

**Primary risks.** Automation bias at speed; checkpoint skipping; error propagation across stages; assumption drift; loss of traceability; "QA theater" mistaken for verification

**Minimum controls.** Define stages with acceptance criteria; require human gates; enforce schema at each stage; maintain prompt logs and version history; use stop conditions; separate drafting vs QA vs approval roles; label all outputs `"Not verified"`

---

The remainder of this section explains how to operationalize this wrapper as a practical "minimum standard" that teams can execute under real conditions. The goal is not perfection. The goal is a defensible baseline: if you follow these controls, you reduce the likelihood of the most severe Level 3 failures and you create enough traceability to identify and correct weaknesses.

### 3.4.1   The "Workflow Run Bundle" (what you must save)

The central artifact of Level 3 governance is the **Workflow Run Bundle**. The bundle is not a bureaucratic attachment; it is the minimum record that proves the work was conducted with discipline. It should be saved for every workflow run, even when the output is not sent externally, because internal artifacts often influence decisions. The bundle is also how the organization learns: it allows you to compare runs, see repeated failure patterns, and refine your templates.

In a governance-first approach, you should treat the bundle as the equivalent of a deal room for a small decision: a compact package that captures what was provided, what was produced, what was uncertain, and what was approved. The bundle must be structured so that a reviewer can audit it quickly. If it is saved as an unstructured set of files with unclear naming, it will not be used. The bundle must therefore be standardized with a predictable directory structure and naming convention (the implementation details appear later in the companion notebook, but the conceptual requirements are here).

A key design choice is to separate **content artifacts** from **governance artifacts**. Content artifacts are the stage outputs: issue maps, hypotheses, risk registers, draft memos. Governance artifacts are the logs and records: intake snapshot, prompt log, assumptions register, verification questions list, reviewer notes, version history. The bundle must include both. A common mistake is to save only the content and assume governance can be reconstructed. It cannot, because many governance decisions are negative space: what was *not* included, what was rejected, what was deferred, what was deemed out of scope. Reviewer notes and version history capture that negative space.

The list below is therefore the minimum deliverable standard. Each item exists because without it a predictable failure mode becomes untraceable or uncorrectable.

---

**Artifact (Save This)**

**Minimum deliverable standard (Workflow Run Bundle).**

1. Intake snapshot (redacted): objective, scope, audience, constraints, facts provided.
2. Workflow identifier and version: template name + version number.
3. Prompt log: per-stage prompts (exact text) and any parameter settings (if applicable).
4. Stage outputs: each stage saved as structured output with *Verification status: Not verified. Human review required.*
5. Assumptions register: consolidated across stages with owners and status.
6. Verification questions list: what must be checked, by whom, and by when.
7. Reviewer notes: approvals, edits, rationale for changes, and unresolved issues.
8. Version history: v0 (model stage outputs), v1 (human edits), v2 (final bundle).

---

It is worth pausing on what "minimum" means here. The bundle is not asking the team to write a full audit report. It is asking the team to preserve the critical evidence that makes the work defensible. Each item contributes a distinct layer of defensibility:

**(1) Intake snapshot.** This establishes what the workflow was asked to do, under what constraints, and based on what facts. Without it, later readers cannot tell whether the workflow answered the right question. The intake snapshot is also the main control against scope creep. If the final deliverable includes content not implied by the intake, the discrepancy becomes visible. Redaction is mandatory because the bundle is saved and often shared. You want the bundle to be safe to circulate internally without exposing unnecessary sensitive details.

**(2) Workflow identifier and version.** This is the control against "template drift." If the workflow changes over time, you must know which version produced which output. Otherwise, you cannot reproduce results or compare runs. Treat workflow templates like internal products: they have versions, they evolve, and changes must be traceable.

**(3) Prompt log.** In agentic workflows, the prompts are effectively the procedure. If you cannot see the prompts used at each stage, you cannot understand how the output was shaped. The prompt log must include the exact text because small changes can have large effects. Parameter settings matter because they influence determinism and style, but even if you do not track parameters, the prompt text is non-negotiable.

**(4) Stage outputs.** Stage outputs are the content artifacts. They must be saved *per stage* rather than overwritten. Overwriting destroys traceability and hides error propagation. Each stage output must include the schema fields (facts provided, assumptions, open questions, stage draft output, verification status) so that reviewers can see what was known and what was guessed at the time.

**(5) Assumptions register.** The assumptions register is the central defense against "assumption drift," where the model introduces premises that later become treated as facts. The register consolidates assumptions across stages and forces ownership. In practice, it should also capture whether an assumption is approved or pending, because pending assumptions should block downstream stages

that rely on them.

**(6) Verification questions list.** This makes verification a tracked activity rather than a vague intention. It also prevents an organizational pathology: sending drafts that depend on unverified claims because "we will check later." The verification list must assign owners and due dates whenever possible. If it cannot assign them, it must at least flag the verification as blocking for external use.

**(7) Reviewer notes.** These are the record of human judgment. They explain what changed, why it changed, and what remains open. Reviewer notes protect the organization from hindsight rewriting: when someone later asks why an assumption was accepted, the record shows the rationale and the decision context.

**(8) Version history.** Versioning enforces authorship accountability. v0 is the model output. v1 is human edits. v2 is the final bundle used for communication. Without version history, teams often "clean up" the model output and then forget what was model-generated and what was human-owned. Version history preserves the distinction.

A common question is whether this bundle is too heavy for routine work. The answer is that the bundle is designed to be lightweight when implemented properly: standardized templates, automated naming, and a habit of saving outputs at the end of each stage. The companion notebook will show how to generate these artifacts with minimal friction. The governance point here is not to increase paperwork; it is to ensure that the organization can defend what it did, learn from what it did, and avoid repeating the same failure patterns.

### 3.4.2   Redaction and minimum-necessary input

Level 3 amplifies the importance of redaction and minimum-necessary input for a simple reason: the workflow produces *more* artifacts and therefore stores *more* text. Every additional stage output is another place where sensitive information can appear, be copied, or be logged. The safest approach is therefore not "be careful," but "design the workflow so that sensitive details are rarely needed." This is exactly the logic of minimum-necessary input: provide only what is required to accomplish the task, and nothing else.

In consulting and corporate strategy work, sensitive information often enters workflows accidentally because teams paste raw materials: client emails, internal spreadsheets, HR lists, deal terms, or draft board materials. These inputs may feel necessary because the model produces better outputs when given more context. But in a governance-first setting, "better" must be balanced against "safe and defensible." Moreover, larger inputs often reduce reviewability: when teams paste large blocks, they cannot easily tell what the model relied on. Smaller, explicit fact lists are easier to review and reduce the chance that a hidden sensitive element is repeated in an output.

At Level 3, redaction is not only about confidentiality; it is also about process discipline. If the workflow is designed to accept facts as structured bullet points, it trains the team to think explicitly about what is known. It reduces the temptation to rely on implicit details embedded in long text. It also makes it easier to detect assumption drift: when facts are a short list, any new claim stands out.

The minimum-necessary principle has several practical implications:

**Use placeholders by default.** Clients, counterparties, products, and internal units should be represented as placeholders (Client A, Segment B, Business Unit C). If the workflow output later needs real names, those names can be inserted in a controlled human editing step, outside the model. This protects the logs and stage outputs from becoming a repository of client-sensitive identifiers.

**Remove unique identifiers.** Emails, phone numbers, account IDs, contract IDs, and any identifiers that can link content to a specific person or transaction should be removed. Level 3 workflows often store outputs for later reuse or training. Unique identifiers create unnecessary risk and are rarely required for strategic reasoning.

**Avoid raw datasets.** Raw financial statements, HR files, customer lists, and proprietary operational data should not be pasted unless explicitly approved and handled in an approved environment. Even when approval exists, governance-first practice encourages summarization into structured facts rather than pasting raw tables.

**Normalize sensitive numbers.** Often, exact numbers are not necessary for drafting a decision memo or structuring a workplan. Ranges (e.g., "mid-teens margin"), indexed values, or qualitative descriptors can be sufficient at early stages. Exact numbers can be inserted later by humans once verified.

**Keep mapping offline.** When placeholders are used, the mapping from placeholders to real names should be stored offline or in a controlled internal system, not in the workflow run bundle. This reduces the chance that the bundle itself becomes sensitive.

Redaction discipline is also a control against another Level 3 failure: the accidental creation of "shadow systems of record." If teams begin to rely on workflow bundles as convenient storage, they may unintentionally create a repository of sensitive information outside official systems. Minimum-necessary input reduces that risk by ensuring that bundles remain governable and safe to store.

The following checklist is a practical minimum. It should be enforced as a habit, not as a one-time training item, because Level 3 workflows make it easy for teams to become sloppy when they are moving quickly.

---

**Checklist**

**Level 3 redaction checklist (minimum).**

1. Use placeholders for client, products, and counterparties (Client A, Business Unit B).
2. Remove unique identifiers (emails, phone numbers, account IDs, contract IDs).
3. Do not paste raw datasets or confidential documents unless explicitly approved.
4. Normalize sensitive numbers (ranges, indexed values) where feasible.
5. Keep a local mapping file (offline) if re-identification is required.

---

Two final governance notes are worth stating explicitly because they are where many teams fail in practice.

First, redaction is not a substitute for judgment about appropriateness. Even if you remove names, a set of deal terms or a description of a restructuring can still be sensitive because the context itself is identifying. Minimum-necessary input therefore includes minimizing distinctive context. For example, instead of pasting "the client is acquiring Company X for $Y billion with a tender offer," you might provide "the client is evaluating an acquisition with a tight timeline and potential regulatory review." The workflow can still draft a risk section and a stakeholder plan based on the general facts, while avoiding details that should remain confined to approved systems.

Second, redaction must be paired with labeling and access control. The Workflow Run Bundle is designed to be saved and shared internally, but not all internal sharing is appropriate. Teams should treat bundles as controlled artifacts. If the organization has a document management system with access controls, bundles should live there. If not, teams should at minimum restrict distribution and maintain a record of who received the bundle. Level 3 is the point where "we saved it somewhere" becomes a governance risk.

The Level 3 wrapper is therefore the smallest set of controls that turns agentic workflows into defensible practice: stage definitions that prevent drift, checkpoints that enforce ownership, logs that preserve traceability, and redaction that keeps the process safe. The workflows that follow in the chapter are designed to sit inside this wrapper. If the wrapper is removed, the workflows will still run, but they will run in the most dangerous way: quickly, confidently, and without evidence.

## 3.5   Core Level 3 workflow patterns (A–D)

Level 3 becomes practical when it is expressed as repeatable workflow patterns that teams can run under real consulting conditions: ambiguity at the start, incomplete facts, multiple stakeholders, and relentless time pressure. The patterns in this section are intentionally minimal. They are not meant to be "the perfect process." They are meant to be a **minimum viable governed workflow** that can be executed consistently, audited later, and improved over time. Each pattern is a small module that can be combined with the others. Together they form a basic Level 3 operating discipline: intake and scope control before analysis, drafting supported by non-verifying QA, assumptions made explicit and gated, and a structured skeptical challenge before anything is treated as decision-ready.

A useful way to read these patterns is to focus on what they prevent. Pattern A prevents scope drift and missing-input hallucination. Pattern B prevents polished drafts from bypassing review. Pattern C prevents assumption drift and error propagation. Pattern D prevents the workflow from becoming a self-confirming machine. These are the four core failure modes that appear repeatedly when organizations attempt "agentic" work without governance-first constraints.

Across all patterns, one requirement stays constant: each stage output must include the same minimum schema fields (facts_provided, assumptions, open_questions, stage_draft_output, verification_status), and each stage must be logged and saved. The workflow is only as defensible as its record.

### 3.5.1 Pattern A: Intake → scope lock → stage plan

Pattern A is the foundation of Level 3 governance because it turns the first moments of a consulting request from informal conversation into a controlled decision about what the workflow is allowed to do. In most organizations, the earliest stage of work is where the most damaging mistakes are made: unclear objectives, mismatched audiences, ambiguous constraints, and hidden sensitivities. AI systems are particularly vulnerable here because they respond to ambiguity by completing it. If you ask a model for "a market entry memo," it will give you one even if the market, the product, and the decision-maker are unspecified. The result can be a beautifully written document that is not merely wrong, but wrong in the way that looks right.

Pattern A blocks that failure by enforcing a simple principle: **no analysis step can begin until the intake and scope lock are approved**. This is a human checkpoint by design. The model can help structure the intake, but it cannot decide what is in scope, what is out of scope, what risks are present, and what facts are required. Those are professional judgments tied to accountability and often to confidentiality.

**Stage 0: Intake form.** The intake stage is a structured capture of the minimal information needed to govern the workflow. It should be short enough to complete under pressure, but complete enough to prevent silent drift. A well-designed intake captures:

- **Objective:** what decision or outcome the deliverable is meant to support.
- **Audience:** who will read it (project team, executive committee, board, client, regulator).
- **Deliverable type:** memo, slide storyline, workshop agenda, stakeholder email, decision record.
- **Constraints:** time, budget, risk appetite, organizational constraints, required tone, non-negotiables.
- **Sensitivity classification:** internal-only vs external; confidentiality level; redaction requirements.
- **Facts provided:** an explicit list of what is known and approved for use.
- **Known unknowns:** the user's own awareness of missing facts or uncertainties.

The governance purpose of Stage 0 is not to gather everything. It is to gather what must be true for the workflow to be legitimate. When teams skip this, they often discover too late that the memo was written for the wrong audience, that a key constraint was ignored, or that sensitive information was included unnecessarily.

**Stage 1: Scope lock (in/out) + risks + required facts.** Scope lock is where the workflow becomes governable. The output of Stage 1 should be a short, explicit statement that includes:

- **In-scope:** what the workflow will produce and what questions it will address.
- **Out-of-scope:** what the workflow will not do (e.g., no external benchmarking, no fact verification, no recommendation).
- **Risk flags:** the types of content that require heightened caution (quantitative claims, regulatory assertions, external commitments).
- **Required facts to proceed:** what must be provided or approved before analysis/drafting

stages can run.

- **Stop conditions:** triggers that will halt the workflow (missing decision-maker, missing constraints, presence of unapproved numbers).

This stage is where many organizations first confront the difference between "AI wrote something" and "we can defend what we did." A scope lock forces clarity. It is also the first place where redaction discipline is reinforced. If the deliverable can be produced with placeholders and general facts, the scope lock should insist on that.

**Stage plan.** Once scope is locked, the workflow proposes a stage plan: a numbered list of stages, each with an objective, required inputs, expected outputs, and a named checkpoint owner. The plan is not merely a project plan; it is the enforcement mechanism. It tells the workflow what it is allowed to do next, and it tells reviewers what to expect. For example, a market entry workflow might include: Stage 2 issue map, Stage 3 alternatives, Stage 4 draft memo, Stage 5 QA pass, Stage 6 red-team critique, Stage 7 human synthesis. A cost transformation workflow might include: baseline taxonomy capture, initiative list, execution risk map, KPI draft, memo draft, QA, human review.

**Human checkpoint: approve scope + facts_provided.** This is the gate. Before proceeding, a human owner must approve the scope lock and the facts_provided list. This is where accountability becomes explicit: the owner is effectively saying, "These are the facts we are using, this is what we are attempting, and this is what we are not attempting." Without this checkpoint, Stage 2 outputs are already compromised. With it, Stage 2 outputs can be challenged against an explicit scope.

Pattern A therefore produces three benefits that compound across the workflow. It prevents the model from inventing context. It prevents the team from accidentally asking the workflow to do something it cannot safely do. And it creates a clear record of what was approved at the start, which is essential when later artifacts are questioned.

### 3.5.2 Pattern B: Draft → QA pass → human review

Pattern B is the minimum viable production line for Level 3 deliverables. It recognizes two realities: (1) models can draft quickly and consistently, and (2) humans cannot review everything with equal attention. Pattern B uses a non-verifying QA pass to improve the efficiency and quality of human review, but it also explicitly guards against the misconception that QA equals verification. The QA pass is a structural and language check, not an evidence check.

**Stage 1: Draft artifact (AI).** The drafting stage must be constrained by the approved facts and the scope lock. It should produce a draft that is clearly labeled as not verified and that includes explicit open questions. The drafting stage is not the place to "fill in missing data." In governance-first practice, missing data is a feature, not a bug. The draft should contain placeholders or "TBD" markers where facts are missing, and it should route those missing items into open_questions and questions_to_verify.

The draft stage should also obey a tone discipline. In consulting, language often implies certainty ("we conclude," "the analysis shows," "it is clear"). In Level 3 drafting, such phrases should be

avoided unless they are grounded in approved facts and verified evidence. The workflow should instruct the model to use cautious language: "based on the facts provided," "preliminary," "subject to verification," "requires confirmation." This is not timid writing; it is accurate writing.

**Stage 2: QA pass (AI as checker; no verification).** The QA pass is an internal audit of the draft as text. It checks whether the draft obeys the workflow's constraints and whether it is structurally safe to hand to a human reviewer. A strong QA pass includes at least four categories of checks:

1. **Schema completeness:** are required fields present; are open_questions populated where needed; is verification_status present.
2. **Facts/assumptions leakage:** does the draft state as fact something that was only an assumption; does it introduce new facts not in facts_provided.
3. **Scope compliance:** does the draft address questions outside the scope lock; does it imply external research or benchmarking.
4. **Language risk:** does the draft imply certainty or verification; does it contain recommendation-like phrasing; does it overclaim.

The QA output should be a short set of findings and required fixes. Importantly, the QA pass must not rewrite the draft automatically in a way that hides what was wrong. It should identify issues, suggest edits, and then either (a) prompt the drafting stage to revise with explicit changes logged, or (b) hand off to the human reviewer with clear "fix these first" notes. The design choice depends on the workflow's maturity, but the governance principle is stable: the QA stage creates transparency about risk, not false reassurance.

**Stage 3: Human review + edits + approval notes.** Human review is the gate that turns drafts into usable work. In Pattern B, the human reviewer's job is not to "polish" language. It is to perform professional skepticism: challenge assumptions, confirm facts, ensure the deliverable is appropriate for the audience, and decide what must be verified before use. The reviewer should also decide whether the workflow run is acceptable to proceed to the next step (e.g., stakeholder circulation) or whether the workflow must halt until verification questions are resolved.

Human review must include recorded notes. This is essential for defensibility. Notes should capture: what was accepted, what was changed, what was rejected, and why. The notes can be brief, but they must exist, because they are the record of human judgment.

Pattern B produces two key outcomes. First, it reduces the chance that a polished draft bypasses skepticism. Second, it makes review scalable by using QA to pre-screen the draft for structural problems and risky language. The workflow becomes safer not because it is clever, but because it forces review at the right place and makes review easier.

### 3.5.3 Pattern C: Assumption register with blocking rules

Pattern C is the control that prevents Level 3 from becoming an error-propagation engine. In multi-step workflows, assumptions introduced early tend to become embedded. If those assumptions

are not made explicit and owned, they will be treated as facts by later stages, by other team members, and by downstream stakeholders. An assumption register is therefore not optional at Level 3; it is the central spine of defensibility.

**What the assumption register is.** The register is a consolidated list of assumptions introduced at any stage of the workflow, with at least four fields: the assumption statement, the owner (human), the stage where it was introduced, and its status (pending/approved/rejected). Many teams also add a risk tag (low/medium/high) based on the impact of the assumption on decisions.

**Why blocking rules are required.** In Level 2, you can list assumptions and still proceed informally. In Level 3, proceeding informally is dangerous because the workflow can move quickly and produce many artifacts that depend on those assumptions. Blocking rules are the mechanism that prevents the workflow from proceeding into high-risk territory without explicit human approval. A blocking rule is a stop condition triggered by certain assumption types.

A practical governance-first blocking approach is to classify assumptions into categories:

- **Low-risk assumptions:** formatting conventions, generic structure choices, non-decision-critical wording. These can proceed with disclosure.
- **Medium-risk assumptions:** directional judgments about feasibility, rough timelines, stakeholder preferences that are plausible but unconfirmed. These should be disclosed and reviewed before external use.
- **High-risk assumptions:** quantitative assumptions (market size, cost savings, ROI), external claims (competitor behavior, regulatory requirements), and commitments (timelines, promises to stakeholders). These must block downstream stages that rely on them until approved or replaced with verified facts.

**How blocking rules work in the workflow.** When a stage introduces an assumption that matches a high-risk category, the workflow must do two things: (1) add the assumption to the register with status "pending," and (2) set stop_conditions_triggered to "yes" for any downstream stage that would rely on that assumption. The workflow then routes the issue to a human checkpoint: approve, revise, or reject. Only after the human decision is recorded does the workflow proceed.

This sounds strict, and it is meant to be. High-risk assumptions are where consulting work becomes decision-dangerous. A workflow that continues past them is essentially recommending through the back door. Blocking rules keep the workflow honest: it cannot manufacture specificity to satisfy the desire for completeness.

**Human checkpoint: approve or revise assumptions.** The checkpoint here is a real decision. The human owner may decide to (a) approve the assumption explicitly as a temporary working premise, (b) revise it to be more conservative, (c) reject it and require verification, or (d) re-scope the deliverable so it does not depend on the assumption. Any of these are valid. What is invalid is proceeding silently.

Pattern C also helps teams manage a practical organizational challenge: different stakeholders have different tolerance for assumptions. Some executives are comfortable with working premises; others demand verified facts. By making assumptions explicit and owned, the workflow allows the

team to tailor communication appropriately. It also reduces later conflict, because the record shows what was assumed and why.

### 3.5.4 Pattern D: Red-team critique before finalization

Pattern D is the skepticism module. It is designed to counteract a structural risk in multi-step workflows: once a workflow has produced a coherent bundle, humans become psychologically inclined to accept it. Coherence creates comfort. Comfort reduces skepticism. Pattern D re-injects skepticism in a structured way by forcing the workflow to generate a critique that highlights weaknesses, missing alternatives, fragile assumptions, and misleading phrasing.

**What red-teaming is at Level 3.** Red-teaming here is not adversarial security testing or "breaking the system." It is a disciplined skeptical review of the deliverable as a decision-support artifact. The model is instructed to take a critical stance: assume the draft could be wrong, incomplete, or misleading, and identify where. Importantly, the red-team stage is itself non-verifying. It cannot claim external truth. Its role is to challenge structure and reasoning, not to introduce new facts.

**What the red-team critique should cover.** A useful red-team critique includes at least five dimensions:

1. **Missing alternatives:** are there plausible options not listed; is the status quo treated explicitly.
2. **Weak assumptions:** which assumptions are doing the most work; what happens if they fail.
3. **Tradeoff imbalance:** does the draft overemphasize benefits or underemphasize risks.
4. **Misleading certainty:** does the language imply verification or inevitability.
5. **Decision readiness:** what must be verified or aligned before the deliverable can support a real decision.

The critique should be specific, not generic. It should point to exact sections where phrasing is risky, where an assumption is unowned, where an open question is missing. In governance-first practice, generic critiques are useless because they do not change behavior. The red-team stage must therefore be structured and tied to the schema.

**Human checkpoint: decide what to incorporate; record rationale.** The red-team critique is not automatically applied. If it were, the workflow would be editing itself in a way that hides the decision. Instead, a human owner reviews the critique and decides what changes to incorporate. This decision is recorded. Sometimes the human will decide that the critique is valid and that the deliverable must be revised. Sometimes the human will decide that the critique reflects a different scope than intended, and will record that rationale. The critical governance point is that the critique creates a moment of forced reflection before the deliverable is treated as final.

Pattern D also has a cultural benefit. It trains teams to expect skepticism as part of the process rather than as personal criticism. When a junior team member critiques a partner's memo, it can become political. When the workflow produces a structured critique, it normalizes challenge. The partner still decides what to do, but the critique becomes part of the workflow discipline.

**How patterns A–D work together.** These patterns are designed to be combined, not chosen. Pattern A prevents the workflow from beginning in ambiguity. Pattern B structures production and review. Pattern C prevents assumptions from silently propagating. Pattern D prevents coherence from becoming complacency. When combined, they create a minimal Level 3 workflow that is safe enough to run repeatedly: intake and scope lock, stage plan, drafting with QA, assumption gating, red-team critique, human synthesis and sign-off.

A final practical note is that Level 3 patterns must be designed for real deadlines. That means the workflow must make it *easier* to do the right thing than to do the wrong thing. If gates are complicated, they will be skipped. If logs are hard to save, they will be lost. If redaction is burdensome, raw content will be pasted. The companion notebook is where the mechanics of automation reduce friction. The chapter's job is to define what must be true for Level 3 governance. These four patterns define that minimum: explicit intake and scope, structured drafting with non-verifying QA, assumptions made explicit with blocking rules, and a forced skeptical critique before finalization.

## 3.6 Mini-cases (Level 3): workflows in real consulting scenarios

Level 3 is where the book's recurring mini-cases become operational rather than illustrative. At Level 1, the same cases are used to demonstrate disciplined drafting. At Level 2, they are used to demonstrate structured reasoning artifacts. At Level 3, they are used to demonstrate something more important: how a consulting team can run a **checkpointed workflow** that produces a defensible **bundle of intermediate artifacts** rather than a single polished document. The emphasis is not "autonomous analysis." The emphasis is **workflow discipline**. The workflow must only use facts provided by humans, must flag missing information explicitly, and must carry a *Verification status: Not verified. Human review required.* posture across every stage output.

In each case below, the purpose is to show what a Level 3 workflow run looks like when executed properly: stage boundaries are clear, gates are enforced, assumptions are consolidated and owned, and handoffs include what must be verified. Each case is presented as a "workflow-run bundle" because that is the unit of defensibility at Level 3. If the team cannot produce the bundle, the work is not Level 3 in the governance-first sense, even if the final memo looks excellent.

A consistent pattern will appear across all cases:

- **Stage 0 (Intake and classification):** capture objective, audience, scope, constraints, sensitivity, facts provided.
- **Stage 1 (Structure):** build the scaffolding (issue map, taxonomy, alternatives framing, stakeholder map).
- **Stage 2 (Reasoning artifacts):** hypotheses, options, tradeoffs, and an explicit evidence/verification plan.
- **Stage 3 (Draft deliverable):** produce the memo, workplan, board pre-read outline, or workshop bundle.

- **Stage 4 (QA pass):** non-verifying checks for structure, consistency, scope compliance, and risky language.
- **Stage 5 (Red-team critique):** skeptical challenge, focused on missing alternatives and fragile assumptions.
- **Human gate(s):** approve scope, approve high-risk assumptions, decide what to incorporate from critique, sign off.

Not every workflow needs all stages, but Level 3 requires at least one explicit gate before drafting and one explicit gate before anything is treated as decision-ready. In practice, teams often add an additional gate after Stage 2, once the key assumptions are visible.

### 3.6.1 Case 1: Market entry (workflow-run decision memo bundle)

A market entry case is an ideal illustration of why Level 3 exists. Market entry requests often begin with vague urgency: "We are considering entering Country X," or "We want to evaluate Segment Y." Teams then rush to produce a memo. The memo looks complete, but it is grounded in implied facts, generic assumptions, and unverified claims. In a governance-first setting, the workflow must resist that impulse and force the team to make uncertainties explicit.

**Objective of the workflow run.** Produce a one-page decision memo draft that frames the market entry question, enumerates structured alternatives, and provides a verification plan, without inventing market data, competitor shares, or regulatory requirements.

**Stage 0: Intake (redacted fact pattern).** A minimal intake might include:

- Objective: "Decide whether to proceed to a full diligence phase for market entry."
- Audience: "Corporate strategy VP and CFO."
- Scope: "Initial screening memo; no external research; identify what must be verified."
- Constraints: "Four-week decision window; risk-averse posture; prefer pilot options."
- Facts provided: "Current product category; current regions served; strategic rationale; high-level capacity constraints."
- Sensitivity: "Internal only; placeholders for target country and partners."

The key control is that the intake explicitly states what is not provided (market size, price points, competitor landscape, regulatory constraints). The workflow must treat these as open questions, not as blanks to fill.

**Human gate (end of Stage 0): approve facts provided and decision question.** The owner confirms that the decision question is correct (screening vs full go/no-go) and that the facts provided are approved.

**Stage 1: Issue map (structure).** The workflow produces an issue map that decomposes the decision into major branches, for example:

- Market attractiveness (demand drivers, willingness to pay, growth)
- Competitive dynamics (incumbents, differentiation, barriers to entry)
- Regulatory and compliance constraints (licenses, data rules, trade restrictions)

- Capability fit (operations, supply chain, talent, partners)
- Economics (unit economics, investment needs, payback)
- Execution and risk (timeline, failure modes, reputational risk)

Crucially, the issue map is not populated with numbers. It is a structure that indicates what evidence would be needed.

**Stage 2: Hypotheses + evidence plan.** The workflow proposes *testable* hypotheses and a verification plan. For example:

- Hypothesis: "Demand is sufficient in the target segment to justify a pilot."
- Evidence needed: "Third-party market reports *or* internal sales proxies *or* targeted customer interviews."
- Open questions: "What is the addressable market under the agreed scope? What is the target price range?"

This stage is where the workflow produces value without inventing facts: it translates uncertainty into a plan.

**Human gate (after Stage 2): approve assumptions or block.** If the workflow introduces high-risk assumptions (e.g., "assume 10% growth"), they must be tagged and blocked. The human owner either removes them, revises them into cautious placeholders, or approves them explicitly as temporary working assumptions.

**Stage 3: Draft memo (decision memo bundle).** The memo draft should follow a stable executive structure:

- Context and objective (based on facts provided)
- Decision needed and timing
- Alternatives (including "do nothing" and "pilot")
- Key considerations and tradeoffs (qualitative unless facts provided)
- Risks and mitigations (explicitly provisional)
- Verification questions (what must be confirmed before a commit decision)
- *Verification status: Not verified. Human review required.* label

**Stage 4: QA pass.** The QA stage flags:

- Any numbers introduced without being in facts provided
- Any language implying verification ("confirmed," "proven," "market is X")
- Any scope creep (e.g., regulatory conclusions stated as fact)
- Missing required fields (assumptions/open questions)

**Stage 5: Red-team critique.** The critique challenges:

- Missing alternatives (e.g., partnership vs direct entry)
- Weakest assumptions (e.g., channel access, timeline feasibility)
- Understated risks (regulatory, reputational, operational constraints)

**Human sign-off.** The reviewer decides what to incorporate and records rationale. The final

bundle includes: the memo draft, the issue map, the hypotheses/evidence plan, QA findings, red-team critique, assumptions register, and verification questions list.

The deliverable is therefore not "a market entry analysis." It is a governed workflow run that produces a decision-support memo draft and a plan to obtain the evidence required for a real decision.

### 3.6.2 Case 2: Cost transformation (workflow-run workplan bundle)

Cost transformation work is particularly vulnerable to AI misuse because it invites benchmark language and invented savings numbers. A model can easily generate a list of initiatives with confident savings estimates, creating the illusion of a feasible plan. Governance-first Level 3 avoids this by treating the output as a workplan scaffold: taxonomy, initiative categories, execution risks, and verification needs.

**Objective of the workflow run.** Produce a structured workplan bundle for a cost transformation initiative, including initiative categories and governance, without inventing baselines or savings.

**Stage 0: Intake.** A minimal intake includes:

- Objective: "Define a transformation workplan and stakeholder kickoff package."
- Audience: "COO, CFO, transformation office."
- Constraints: "No headcount actions stated without approval; union-sensitive; timeline within two quarters."
- Facts provided: "High-level cost areas; organizational structure; known constraints (e.g., systems limitations)."

**Human gate: approve constraints and sensitive boundaries.** This is critical: many cost programs are politically and reputationally sensitive. The scope lock must explicitly forbid certain content (e.g., recommending layoffs) unless the human owner approves.

**Stage 1: Cost taxonomy framing.** The workflow produces a structured taxonomy suitable for the organization (SG&A, COGS, overhead, procurement, IT, facilities, etc.), but it does not populate amounts. It identifies what baseline data is required to quantify opportunities.

**Stage 2: Initiative list (non-recommending).** The workflow proposes initiative *categories* and examples, framed as options for assessment:

- Procurement: supplier rationalization, contract renegotiation
- Operations: throughput improvements, process standardization
- IT: application rationalization, license optimization
- Org design: spans and layers review (without prescriptive headcount actions)

Each initiative includes: required data, typical risks, dependencies, and verification questions. The workflow explicitly avoids savings estimates unless provided by the human.

**Stage 3: Risk register + governance.** This stage produces execution risks (change fatigue, data quality, stakeholder resistance), governance risks (scope creep, unclear decision rights), and

control mechanisms (weekly steering, decision logs, benefit tracking discipline).

**Stage 4: Stakeholder email / kickoff note.** The workflow drafts a cautious kickoff message that emphasizes the process, the need for validated baselines, and the non-commitment posture.

**Stage 5: QA.** The QA pass flags:

- Any implied quantified savings
- Any prescriptive HR actions not in scope
- Overconfident language about outcomes

**Stage 6: Human review and approval notes.** Human reviewers adjust tone, enforce sensitivity boundaries, and confirm which initiatives are appropriate to propose for assessment.

The final bundle is a workplan scaffold that can be executed: taxonomy, initiative categories with evidence needs, risk register, governance plan, and kickoff comms. It is useful without inventing financial impact.

### 3.6.3   Case 3: Capital allocation (workflow-run board pre-read bundle)

Capital allocation cases are governance-intensive because they often involve large irreversible commitments and board-level scrutiny. They also tempt the model to produce quantitative narratives: IRR, payback, NPV, and "optimal" portfolio choices. In a governance-first Level 3 workflow, the model can structure the alternatives and the decision logic, but it cannot compute or assert financial outputs unless the numbers are provided and verified.

**Objective of the workflow run.** Produce a board pre-read structure that frames alternatives, tradeoffs, and verification needs, without making the decision or inventing financial metrics.

**Stage 0: Intake.** Minimal intake includes:

- Objective: "Prepare a board pre-read outline for a capital allocation decision."
- Audience: "Board of directors; CFO; investment committee."
- Constraints: "No recommendation; present options and decision criteria; conservative risk posture."
- Facts provided: "List of candidate uses of capital; strategic priorities; any known constraints (debt covenants, liquidity targets)."

**Human gate: approve decision criteria and constraints.** The board pre-read must reflect the organization's actual decision criteria (strategic fit, risk-adjusted return, resilience, optionality), not generic textbook criteria.

**Stage 1: Alternatives framing.** The workflow enumerates alternatives explicitly, including:

- Reinforce core operations (capex, productivity)
- Strategic acquisition
- Organic growth investment
- Debt reduction
- Share repurchase or dividend
- Liquidity reserve (do nothing now)

The key control is that "status quo" is treated as a real alternative.

**Stage 2: Tradeoff mapping and scenario template.** The workflow builds a tradeoff matrix that compares alternatives along agreed criteria, and produces a scenario template that shows what would need to be true for each alternative to be favorable. This is where the workflow creates value: it frames uncertainty, it does not solve it.

**Stage 3: Risks/assumptions register.** Capital allocation is assumption-heavy (growth outlook, cost of capital, execution risk). The workflow consolidates assumptions and flags high-risk ones for blocking. Any numeric assumption must be either provided or left as an open question.

**Stage 4: Draft board pre-read outline.** The output is a structured outline with placeholders:

- Decision statement and timing
- Alternatives and strategic rationale
- Decision criteria and weighting (if provided)
- Risk and resilience considerations
- Information required before final decision (verification plan)
- Appendices (models to be provided by finance)

**Stage 5: Red-team critique.** The critique challenges whether alternatives are complete, whether the decision criteria are biased, and whether risks are understated.

**Human synthesis.** The human owner decides how to present the material to the board and ensures that the pre-read does not imply verified financial results that are not yet validated.

The bundle becomes a governance tool: it frames a board decision with clarity and discipline, rather than providing a model-generated "optimal answer."

### 3.6.4 Case 4: Operating model redesign (workflow-run workshop bundle)

Operating model redesign is a classic consulting scenario: centralization vs decentralization, governance tradeoffs, decision rights, and change management. It is also a scenario where models can produce plausible RACI charts and org structures that look credible but may be politically or operationally infeasible. Level 3 governance focuses on producing a workshop bundle that supports human discovery and alignment, not a model-imposed design.

**Objective of the workflow run.** Produce a workshop bundle (agenda, interview guide, provisional RACI draft, risk register) to support operating model redesign discussions.

**Stage 0: Intake.** Minimal intake includes:

- Objective: "Run a workshop to align on target operating model principles and decision rights."
- Audience: "Functional leaders and transformation team."
- Constraints: "No final org design; workshop is for alignment; avoid naming individuals; place-holders only."
- Facts provided: "High-level current operating issues; known pain points; any non-negotiable constraints."

**Human gate: approve sensitivity constraints.** Operating model work often touches

personal and political issues. The workflow must explicitly forbid identifying individuals and must avoid prescriptive org charts.

**Stage 1: Stakeholder map.** The workflow produces a stakeholder map (roles, interests, concerns) at a high level, with placeholders.

**Stage 2: Workshop agenda.** The agenda is structured to elicit decisions, not to present a model's design:

- Current-state pain points and objectives
- Design principles (what must be true)
- Decision rights mapping (who decides what)
- Options: centralize vs federate vs hybrid
- Governance and escalation paths
- Next steps and evidence needed

**Stage 3: Interview guide.** The workflow produces an interview guide that supports discovery: questions about bottlenecks, accountability, cross-functional friction, and failure modes. It avoids leading questions that push toward a preconceived solution.

**Stage 4: Provisional RACI draft (clearly labeled).** The workflow drafts a provisional RACI based on generic roles, explicitly labeled as *Verification status: Not verified. Human review required.* and "for discussion only." This is an important governance posture: the RACI is not a recommendation; it is a discussion scaffold.

**Stage 5: Risk register.** The workflow produces risks: change resistance, role confusion, governance overload, accountability gaps, and transition risks.

**Stage 6: QA and human review.** QA flags any language that implies final decisions or that names individuals. Human review adjusts the bundle for organizational nuance.

The final bundle is a controlled enablement package for a workshop. It accelerates preparation, but it does not outsource design decisions. It also provides a traceable record: what assumptions were used to draft the provisional RACI, what open questions remain, and what must be verified through interviews and data.

**What these cases teach at Level 3.** Across all four scenarios, the common lesson is that Level 3 is not about asking the model to do more "analysis." It is about asking the organization to do more governance. The workflows produce bundles that are useful precisely because they make uncertainty explicit, separate facts from assumptions, and force checkpoints before anything is treated as decision-ready. When run properly, these workflows increase speed where speed is safe (drafting structures, generating checklists, drafting agendas) and slow down where speed is dangerous (approving scope, approving assumptions, committing to conclusions). That is the defining discipline of Level 3.

## 3.7 Risks and controls taxonomy (Level 3)

Level 3 is the point where AI capability begins to look operational. A workflow can produce a chain of artifacts quickly, in formats that resemble what a consulting team would produce under normal conditions: intake summary, issue map, hypotheses, draft memo, QA notes, red-team critique. This is exactly why Level 3 needs a sharper risk and controls taxonomy than Levels 1 and 2. The most serious failures at this level are not "the model wrote something odd." They are failures of *process integrity.* A workflow that produces polished artifacts can create a powerful illusion: that steps were followed, that diligence occurred, that reasoning was tested. When that illusion is not backed by real governance, the workflow becomes a credibility generator rather than a decision-support mechanism.

A taxonomy is useful only if it is actionable. The risk list below is therefore written as "failure modes you should expect," not "risks you might imagine." These are predictable breakdowns that appear in real teams when they attempt to scale AI usage from single outputs to workflow runs. The control set that follows is the minimum enforceable countermeasure. It is not a mature enterprise control framework; it is a baseline that makes Level 3 defendable in practice.

### 3.7.1 Failure modes you should expect

1. Checkpoint erosion (steps skipped under deadline pressure).
2. Error propagation (one flawed assumption infects multiple downstream artifacts).
3. QA theater (consistency checks mistaken for verification of truth).
4. Workflow drift (the "agent" starts answering beyond the defined scope).
5. Recordkeeping gaps (missing prompts, missing stage outputs, unclear versions).

These five failure modes are related. They rarely occur alone. Most Level 3 incidents are a cascade: checkpoint erosion allows a workflow to proceed without approval; a weak assumption slips through; error propagation embeds that assumption into multiple artifacts; QA theater creates false reassurance; and recordkeeping gaps prevent the team from reconstructing what happened. The outcome is not only an incorrect artifact, but a loss of organizational confidence in the process. To prevent that cascade, it helps to understand each failure mode as a pattern of behavior and incentives.

**1. Checkpoint erosion.** Checkpoints fail because they are friction, and consulting environments are optimized for speed. When a deadline approaches, teams rationalize skipping a gate: "We just need a draft," "We will review later," "It is internal." The workflow then proceeds as if approval occurred. Over time, skipping becomes normalized. Once normalized, the workflow is no longer governed; it is autonomous in practice, even if not in design.

Checkpoint erosion is not a moral failure; it is a design failure. If checkpoints are vague, they will be skipped. If checkpoints are difficult to execute, they will be skipped. If checkpoint clearing is not recorded, it will be assumed. Therefore, governance-first Level 3 treats checkpoints as engineered controls: they must be explicit, simple, and recorded. A checkpoint is not "someone should review."

It is a stop condition until a named person approves.

There is a second form of checkpoint erosion that is more subtle: *rubber-stamp approval.* A reviewer may "approve" without actually reading because the workflow output looks professional. This is an instance of automation bias. It can be mitigated by requiring reviewers to confirm specific elements: "I confirm facts provided are accurate," "I confirm assumptions are acceptable," "I confirm open questions are listed," "I confirm no verification is implied." In other words, the checkpoint must be structured to demand thought, not just a signature.

**2. Error propagation.** Multi-step workflows turn small errors into large ones. In a single memo, an invented claim might be corrected in a human edit. In a workflow run, the invented claim can appear in an issue map, then be reflected in hypotheses, then be embedded in the narrative, then be repeated in a risk register, and finally be referenced in a stakeholder email. By the time the human reviewer sees the final draft, the claim appears "supported" because it appears everywhere. The reviewer may not notice that it originated as an unowned assumption.

Error propagation is particularly dangerous when the error is not obviously "wrong," but plausibly generic: a typical timeline, a typical market growth rate, a plausible cost savings range. These are exactly the kinds of claims models produce fluently. Because they are plausible, they are often not challenged, especially when teams are under time pressure. This is why the assumption register and schema enforcement are non-negotiable at Level 3. If every stage output must explicitly list assumptions and open questions, then propagation becomes visible. If assumptions are consolidated across stages, then reviewers can inspect the "assumption spine" of the bundle before trusting any downstream artifact.

Error propagation also interacts with scope drift. A workflow that begins with a narrow scope can gradually incorporate new sub-questions. Each expansion introduces new assumptions. Without explicit scope lock and stop conditions, the workflow can "solve" questions it was never authorized to address, and then propagate those answers downstream as if they were part of the original mission.

**3. QA theater.** QA theater is when consistency checks are mistaken for verification of truth. At Level 3, teams often introduce an AI QA pass because it is easy and because it improves surface quality. The QA pass checks that the memo is internally consistent, that it includes required sections, that it does not contradict itself, and that it uses appropriate tone. These checks are valuable. The failure occurs when the team begins to treat the presence of a QA report as evidence that the content is "validated."

This is a predictable cognitive trap. Humans treat process artifacts as signals of safety. When a workflow produces a "QA passed" message, it feels like a green light. But internal consistency is not truth. A perfectly consistent memo can be perfectly wrong. A contradiction scan cannot detect missing evidence. A completeness check cannot detect a false premise. Therefore, Level 3 governance must enforce the language of QA. QA outputs must explicitly state: "This is not verification. This is structural review only." QA should also highlight where verification is required, not only where text is inconsistent.

The most dangerous form of QA theater is when teams allow QA to substitute for skeptical

human review. Reviewers may skim the QA report and assume the output is safe. Governance-first Level 3 therefore designs QA as a helper for reviewers, not as a gatekeeper. The gatekeeper remains the human checkpoint, and the checkpoint must require explicit confirmation that verification has not occurred.

**4. Workflow drift.** Workflow drift is when the "agent" starts answering beyond the defined scope. Drift can be driven by the model (it attempts to be helpful) or by the user (they ask additional questions mid-run). Drift is dangerous because it undermines accountability: the workflow run no longer corresponds to the approved intake. It also undermines defensibility: the output may contain claims that were not intended to be made, especially claims that imply recommendation or verification.

Workflow drift often presents as incremental expansion: the draft memo begins to include "industry benchmarks," the risks section begins to include "regulatory requirements," the alternatives section begins to rank options. None of these expansions may be explicitly requested, but the model includes them because they are common in similar memos. In consulting, this is especially risky because conventional memo formats often assume access to data and research. A model trained on such formats will imitate them.

Governance-first controls against drift rely on scope lock, schema enforcement, and QA flags. Scope lock sets the boundary. Schema enforcement forces the model to label any new claims as assumptions or open questions. QA checks for scope creep and for recommendation-like phrasing. Stop conditions trigger when the model attempts to cross the boundary: for example, if the model introduces numeric benchmarks not present in facts_provided, the workflow halts and escalates.

**5. Recordkeeping gaps.** Recordkeeping gaps are the failure mode that turns a manageable mistake into an organizational incident. Even well-run workflows will produce imperfect drafts. The organization can learn from them only if it can reconstruct what happened. Missing prompts, missing stage outputs, overwritten files, unclear version labeling, and absent reviewer notes destroy reconstruction. When reconstruction is impossible, the organization either overreacts ("we cannot use this at all") or underreacts ("we cannot prove it was wrong"), both of which harm governance.

Recordkeeping gaps are common because teams focus on the final deliverable and treat intermediate artifacts as disposable. Level 3 reverses that assumption: intermediate artifacts are part of the record because they show how the deliverable was produced. Prompt logs are part of the record because they define the procedure. Version history is part of the record because it shows what was model-generated versus human-edited. Reviewer notes are part of the record because they document judgment. Without these, Level 3 is not defensible.

Recordkeeping gaps also interact with confidentiality risk. If teams store outputs inconsistently, they may copy them into uncontrolled locations. A governed bundle reduces that risk by standardizing where artifacts live and how they are labeled.

### 3.7.2 Minimum control set (practical and enforceable)

The minimum control set is the smallest set of controls that addresses the five failure modes above. Each control is enforceable without heavy infrastructure. The controls are designed to be used in every workflow run, not only in "big" engagements. If the controls are optional, they will be ignored precisely when they are most needed: when the team is under pressure.

---

**Checklist**

**Minimum controls (use in every Level 3 workflow run).**

1. **Stage schema requirement:** facts_provided / assumptions / open_questions / stage_draft_output / verification_status.
2. **Human checkpoints:** defined gates that block downstream stages until approved.
3. **Stop conditions:** if missing required inputs or high-risk claims appear, halt and escalate.
4. **Separation of duties:** drafting, QA, and approval roles must be named (even if same person in small teams).
5. **Audit trail:** prompt log + stage outputs + reviewer notes + version history saved as a bundle.

---

The checklist is compact, but each item carries important design implications.

**1. Stage schema requirement.** This is the primary defense against hidden assumptions and false completeness. The schema forces the workflow to declare what is known, what is assumed, what is unknown, and what has been drafted. The schema also prevents accidental verification language by making "verification_status" explicit. If every stage output must include verification_status, then the workflow cannot easily pretend that a later stage somehow "verified" earlier claims. The schema requirement also makes review efficient: reviewers know exactly where to look for assumptions and open questions.

A practical governance note: schema compliance should be treated as a hard acceptance criterion. If a stage output lacks assumptions or open questions, that is not "nice to have." It is a failure that should trigger re-run or revision. Otherwise, the schema becomes decorative.

**2. Human checkpoints.** Human checkpoints are the control against checkpoint erosion and against de facto autonomy. The workflow must define at least two checkpoints: one after intake/scope lock and one before externalization or finalization. Many workflows add a checkpoint after hypotheses/evidence planning, because that is where high-risk assumptions often appear. The key point is that checkpoints are not passive. A checkpoint must block downstream stages until approved. This must be true even when the same person is both drafter and approver in a small team. In that case, the checkpoint still exists as a named role and a recorded decision, because the goal is not to add people, but to add accountability.

**3. Stop conditions.** Stop conditions are what make a workflow safe under pressure. They create an engineered refusal to proceed when proceeding would manufacture certainty. Stop conditions should trigger in at least three situations:

- Missing required inputs (e.g., no decision-maker, no scope lock, no constraints).
- High-risk claims appear (numbers, external claims, commitments) without being in facts_provided.
- Scope is exceeded (the workflow begins to answer questions not approved in intake).

When a stop condition triggers, the workflow must escalate to a human owner, and the record must show that escalation. This prevents the "bad process at speed" failure where the workflow produces a finished-looking bundle from incomplete inputs.

**4. Separation of duties.** This control reduces automation bias and rubber-stamp approval. It is a cultural and procedural control: the workflow names roles for drafting, QA checking, and approval. In mature teams, these are different people. In small teams, they may be the same person, but the separation still matters because it forces the person to switch modes: first produce, then critique, then approve. That mode switch is a practical defense against complacency. It also creates a record: if the same person performed all roles, that fact is visible, and the organization can evaluate whether additional review is warranted for high-risk outputs.

**5. Audit trail bundle.** The audit trail is the ultimate backstop. It addresses recordkeeping gaps directly. The bundle must include the prompt log, stage outputs, reviewer notes, and version history. This is what allows reconstruction, learning, and defensibility. Without it, Level 3 cannot scale because the organization cannot prove that governance occurred. The audit trail also enables controlled improvement: workflows can be refined based on observed failures rather than on anecdotes.

A final point is that minimum controls are only "minimum" if they are reliably executed. The biggest governance risk at Level 3 is not the absence of controls in policy. It is the presence of controls that are not followed. Therefore, the controls must be designed to be easy to comply with. The companion notebook operationalizes this by producing structured outputs and logs by default. But the conceptual commitment must exist in the chapter: Level 3 governance is a discipline of *execution*, not of intention.

If a team can reliably apply these five controls, it will prevent most Level 3 failures from becoming serious incidents. The workflow may still produce imperfect drafts, but the organization will (1) know where uncertainty exists, (2) know where assumptions were introduced, (3) know who approved progression, and (4) have a record that enables verification and correction. That is the defensible baseline for agentic workflows in consulting and corporate strategy work.

## 3.8 Prompt patterns and exercises (copy/paste)

At Level 3, prompt patterns are not "clever instructions." They are **procedural building blocks**. The reader should treat them as standardized operating procedures for governed workflow runs. The goal is not to extract brilliance from the model; the goal is to enforce consistency, prevent drift, and make review feasible. This section therefore provides copy/paste templates that implement Level 3 controls directly: intake classification and routing, stage execution with stop conditions, and a non-verifying QA pass. The templates are intentionally conservative. They assume the model

is persuasive and fallible, and they assume the team is busy. If a template allows ambiguity, the model will fill it. If a template allows silent progression, the workflow will drift. If a template allows "helpful" additions, the workflow will invent. Therefore, the templates aim for mechanical clarity.

A second purpose of this section is cultural. In many consulting teams, prompting becomes an informal craft practiced by whoever is most comfortable with the tool. That is a recipe for inconsistent governance. Level 3 requires a shared internal discipline: the team should have a small library of approved prompts that match approved workflows. This is why the templates embed schema requirements, verification posture, and stop conditions explicitly. The prompts are designed to be used verbatim, not reinterpreted each time.

There is also an important boundary: these prompts are designed for **workflow execution**, not for autonomous tool use or independent research. If a team wants external research, it must be introduced as a separate controlled process with human verification, and it must be reflected as facts_provided. The model may assist in structuring verification questions, but it must not fabricate sources or numbers. If you adopt these prompts, you are adopting a discipline: the model outputs are drafts, and humans own the truth.

### 3.8.1 Prompt Template 1: Workflow intake + routing (Stage 0)

Stage 0 prompts are about classification and governance setup. They are the front door of the workflow. A strong Stage 0 output does not attempt to answer the problem. It produces: (1) a clean fact list, (2) explicit missing inputs, (3) risk flags, and (4) a stage plan with checkpoints and owners. In practice, Stage 0 is where you prevent most downstream failures. If the intake is ambiguous, everything after it will be fragile. If the intake is strict, later stages can be efficient.

A practical recommendation is to treat Stage 0 as a required ritual. Even if the request seems simple, run Stage 0 first. The time you spend here is recovered later in fewer revisions and fewer stakeholder surprises. Stage 0 is also where redaction discipline is reinforced: "paste redacted context" is not a suggestion; it is a control.

### 3.8.2 Prompt Template 2: Stage output with stop conditions (Stage N)

Stage execution prompts must enforce the core Level 3 rule: the workflow must be allowed to stop. Stop conditions prevent the model from "helping" by inventing. They also prevent the team from pushing the workflow into decision territory when inputs are missing. The key design feature is that the output format includes `stop_conditions_triggered` and that the prompt instructs the model to set it to "yes" when required inputs are missing. This creates a mechanical mechanism for gating.

Another key feature is the separation between **approved inputs** and **prior stage context**. Approved inputs are the only facts that can be treated as true. Prior stage context is summary-only and may contain assumptions or open questions. This separation reduces the chance that the model treats earlier drafts as verified facts. In governance-first practice, the workflow must be explicit: a prior stage output is not evidence; it is a draft.

Stage prompts also require that the model runs basic QA checks on its own output (consistency, completeness, contradiction scans). This is helpful, but it must be framed as non-verifying. The model can check whether it followed the schema and whether it contradicted itself. It cannot check truth.

### 3.8.3 Prompt Template 3: QA pass (non-verifying)

The QA template is the "skeptic within the workflow," but it is not a verifier. Its job is to protect the team from the most common rhetorical and structural dangers: missing fields, contradiction between facts and assumptions, scope creep, and language that implies certainty. The QA template should be run consistently. The output should be short and actionable: findings, fixes, unsafe phrases. In practice, reviewers should not accept a draft that has not passed through QA, because QA is one of the few controls that can catch problems before humans invest time reading.

However, teams must not treat QA as approval. QA is a filter, not a gate. The gate remains the human checkpoint. In fact, QA outputs should explicitly remind readers of that. It is useful for the QA output to include a "verification required" reminder even if the prompt already specifies it, because redundancy is a safety feature under pressure.

### 3.8.4 Exercise set (team training)

Exercises are not optional in Level 3 adoption because the main failure mode is not misunderstanding; it is drift under pressure. Teams understand, intellectually, that they should not skip checkpoints. Then the deadline arrives, and they do. The only reliable countermeasure is rehearsal. Exercises should therefore simulate real conditions: ambiguous inputs, incomplete facts, stakeholders pushing for certainty, and time pressure. The exercises below are designed to be run in short sessions (30–60 minutes each) with explicit artifact saving. Every exercise requires saving the Workflow Run Bundle, because the goal is to train process, not output.

A useful way to run these exercises is in pairs: one person plays "requester under pressure" and tries to push the workflow past gates; the other plays "workflow owner" and must enforce stop conditions and checkpoint discipline. This mirrors real engagement dynamics.

**Exercise 1: Market entry intake stress test (Stage 0 discipline).** Provide a deliberately sparse, redacted request: "We want to enter Market A in six months; draft a decision memo for the CFO." Run Stage 0 only. The team must produce: facts_provided, open_questions, risk_flags, and a stage plan. The exercise is successful only if open_questions is substantial and if the stage plan includes a scope lock checkpoint. The team should then critique the Stage 0 output: did it incorrectly assume market size, customer segments, or regulatory requirements? Did it correctly refuse to guess? Save the Stage 0 output as part of the bundle.

**Exercise 2: Draft with stop condition triggered (Stage N refusal).** Using the same market entry case, attempt to run the drafting stage without providing key facts (e.g., no audience nuance, no alternatives, no constraints). The stage prompt should cause `stop_conditions_triggered` to be "yes." The exercise is successful only if the model refuses to produce a complete-looking memo

and instead lists missing inputs. The team should then provide the missing facts and re-run the stage, producing `stop_conditions_triggered:  "no"`. Save both runs (the refusal run and the compliant run) and compare. This trains the team to treat refusal as success, not as failure.

**Exercise 3: Assumption register with blocking rules (propagation control).** Provide a cost transformation request and include one tempting but unverified claim, such as "we believe SG&A is bloated." Run the workflow through initiative listing and risk register stages. The exercise is successful only if the workflow identifies that "SG&A is bloated" is an assumption and flags what evidence would be needed (baseline breakdown, benchmarking, internal targets). Then introduce an explicit high-risk assumption: "assume $50M savings." The workflow must block downstream drafting steps unless a human approves. The team should practice approving the assumption as "temporary working premise" and also rejecting it as "requires verification." Save the assumptions register and the checkpoint note.

**Exercise 4: QA theater detection (QA is not verification).** Provide a draft memo (intentionally flawed) that contains an invented benchmark and a phrase like "we confirmed." Run the QA prompt. The exercise is successful only if the QA pass flags the invented benchmark as unsupported, identifies "confirmed" as unsafe phrasing, and lists required fixes. Then test the team: ask a participant to argue that QA means it is safe. Another participant must rebut, citing the governance posture: QA checks structure, not truth. This exercise is cultural; it trains language discipline and prevents false reassurance.

**Exercise 5: Workflow drift and scope creep (boundary enforcement).** Provide a capital allocation request scoped to "outline only." Then, mid-workflow, ask the model for "the best option" or "expected IRR." The workflow owner must enforce the scope lock: the request is out of scope and should be handled as an open question or as a separate verified analysis process. The exercise is successful only if the workflow refuses to recommend, reframes the request into alternatives and verification needs, and records the attempted scope expansion in reviewer notes. Save the bundle, including the note that scope creep was attempted and rejected.

**Exercise 6: Checkpoint stress test (deadline pressure simulation).** Simulate a realistic pressure scenario: "The CEO meeting is in two hours; we need the memo now." Run a mini workflow (intake, scope lock, draft, QA) and intentionally attempt to skip the scope lock checkpoint. The workflow owner must refuse to proceed. The exercise is successful only if the gate holds. Then allow the checkpoint to be cleared quickly with explicit approval notes and proceed. The point is to train the team to make gating fast, not to remove gating.

Across all exercises, require that the team saves a Workflow Run Bundle. After each exercise, conduct a brief retrospective using the bundle:

- Where did assumptions appear?
- Where did the model attempt to imply certainty?
- Where did the workflow drift?
- Did stop conditions trigger appropriately?
- Are prompts, outputs, and versions clearly saved?

If the team can run these exercises and produce bundles that clearly show what was done, what was assumed, what was not verified, and who approved what, then the team is practicing Level 3 correctly. If the team produces polished drafts without visible gates and records, then the team is not practicing Level 3, regardless of output quality.

The templates below are therefore more than copy/paste prompts. They are the minimal operating system of Level 3: they standardize intake, enforce stop conditions, separate drafting from QA, and train teams to treat governance as part of the work, not as an afterthought.

### 3.8.5   Prompt Template 1: Workflow intake + routing (Stage 0)

**Prompt (Copy/Paste)**

```
ROLE: You are a workflow intake assistant for consulting/strategy work.
Do not recommend. Do not verify. Do not invent facts, data, benchmarks, or citations.


CONTEXT (provided by user):
[PASTE REDACTED CONTEXT HERE]


TASK:
Classify the request and propose a stage plan with mandatory human checkpoints.
If required facts are missing, list them as open_questions (do not guess).


OUTPUT FORMAT (STRICT):
facts_provided: [bullet list]
assumptions: [bullet list]  (leave empty unless explicitly allowed by user)
open_questions: [bullet list]
risk_flags: [bullet list]   (e.g., quantitative claims, external commitments, regulatory claims)
stage_plan: [numbered list of stages with objective + required inputs + checkpoint owner]
verification_status: "Not verified"
questions_to_verify: [bullet list]
```

### 3.8.6 Prompt Template 2: Stage output with stop conditions (Stage N)

**Prompt (Copy/Paste)**

```
ROLE: You are executing Stage [N] of a governed workflow.
You must not invent facts. You must stop if required inputs are missing.
You must separate facts from assumptions and label output as not verified.


APPROVED INPUTS:
facts_provided:
[PASTE]


PRIOR STAGE CONTEXT (summary only):
[PASTE]


STAGE OBJECTIVE:
[PASTE]


STOP CONDITIONS (must obey):
- If any required input is missing, output stop_conditions_triggered: "yes" and list what is missing.
- If you detect high-risk claims (numbers, external facts, commitments), flag them as open_questions.
- Do not proceed past the stage objective.


OUTPUT FORMAT (STRICT):
facts_provided: [bullet list]
assumptions: [bullet list]
open_questions: [bullet list]
stage_draft_output: [the stage artifact]
stop_conditions_triggered: ["yes" or "no"]
qa_checks_run: [bullet list] (consistency/completeness/contradictions only)
verification_status: "Not verified"
questions_to_verify: [bullet list]
```

### 3.8.7 Prompt Template 3: QA pass (non-verifying)

```
Prompt (Copy/Paste)

ROLE: You are a QA checker. You do not verify truth. You only check structure and consistency.
Do not add new facts.

INPUT (stage output):
[PASTE stage_draft_output and its facts/assumptions/open_questions]

TASK:
Identify:
1) missing required fields,
2) contradictions between facts and assumptions,
3) places where the draft implies verification or certainty,
4) scope creep beyond the stated objective.

OUTPUT FORMAT:
qa_findings: [bullet list]
required_fixes: [bullet list]
unsafe_phrases_to_remove: [bullet list]
verification_status: "Not verified"
```

### 3.8.8 Exercise set (team training)

ckpoint stress test": simulate deadline pressure and ensure the gate still holds.

## 3.9 Conclusion and transition to Level 4 (Innovators)

### 3.9.1 Summary of main takeaways

Level 3 is where AI use in consulting and corporate strategy becomes operational. The organization stops treating the model as a drafting convenience and starts treating it as a workflow operator that executes a predefined sequence of stages. This shift is not primarily about power. It is about **discipline**. A Level 3 team is not defined by how sophisticated its prompts are, or by how polished its outputs look. It is defined by whether it can run a governed workflow repeatedly, under real pressure, and still produce outputs that are traceable, reviewable, and explicitly *Verification status: Not verified. Human review required..*

The core lesson is that speed changes the risk profile. At Level 1, a flawed paragraph is often caught in editing. At Level 2, flawed reasoning is more visible because structure exposes assumptions. At Level 3, the workflow can produce many artifacts quickly, and that speed can create an illusion of diligence: the bundle looks complete, so people treat it as correct. The primary Level 3 risks therefore are not cosmetic errors. They are failures of process integrity: checkpoint erosion, error

propagation across stages, QA theater mistaken for verification, workflow drift beyond scope, and recordkeeping gaps that prevent reconstruction. Each of these failure modes is predictable, which is why Level 3 governance is not optional. It is the condition of safe use.

The chapter's governance wrapper is designed to be minimal but enforceable. It requires stage definitions with acceptance criteria, mandatory human gates, schema enforcement at every stage, persistent logs and version history, explicit stop conditions, and named roles for drafting, QA, and approval. These controls transform the workflow from "multi-step prompting" into a defendable process. They also create a practical operational artifact: the Workflow Run Bundle. The bundle is the unit of accountability at Level 3. It captures what facts were provided, what assumptions were introduced, what open questions remain, what the model drafted at each stage, what QA flagged, what red-team critique challenged, and what humans approved or changed. Without the bundle, a Level 3 workflow is not auditable and therefore not governable.

A second lesson is that governance-first Level 3 is intentionally anti-autonomous. The model is allowed to stop. In fact, the workflow is designed to stop when it should: when required inputs are missing, when high-risk claims appear, when scope creep is detected, or when assumptions are unowned. This is where teams often need a mindset adjustment. Many users interpret a workflow "refusal" as failure. At Level 3, refusal is often success, because it prevents the organization from manufacturing certainty out of absence. The workflow's job is not to always produce a finished-looking artifact; it is to produce a truthful representation of what is known, what is assumed, and what must be verified before decisions are made.

Finally, Level 3 makes explicit the professional boundary that must remain non-negotiable: AI does not decide, recommend, or verify. It can structure, draft, and check for consistency, but it cannot substitute for judgment. The chapter's patterns (intake to scope lock, draft to QA to human review, assumption register with blocking rules, red-team critique before finalization) are all variations of the same principle: **keep humans responsible, and make the workflow produce evidence of that responsibility**. If Level 2 made thinking inspectable, Level 3 makes execution inspectable. That is the maturity gain.

### 3.9.2 What comes next (preview of Level 4)

If Level 3 is workflow discipline, Level 4 is the beginning of **institutional capability**. At Level 4, teams stop building workflows ad hoc and start creating reusable assets that can be deployed consistently across engagements, workstreams, and teams. This is where the organization begins to treat AI-enabled workflows as internal products: standardized templates, controlled variants, evaluation harnesses, release notes, and supervision frameworks. The shift is subtle but consequential. A Level 3 workflow run is a governed instance. A Level 4 asset is a governed *reusable artifact* that shapes many future runs.

This is why Level 4 governance must expand. Once workflows become reusable, the risk is no longer limited to one run. A design flaw in a template can propagate across many engagements. A misleading prompt pattern can be copied and reused as "best practice." A weak QA stage can

become standard operating procedure. Reuse scales impact, and scaled impact demands stronger controls. Level 4 therefore introduces change management as a first-class requirement: workflow templates need versioning, change logs, and controlled release. When a template changes, the organization must know what changed and why, and it must be able to assess whether the change improved reliability or introduced new failure modes.

Level 4 also introduces evaluation as a routine activity rather than an afterthought. At Level 3, teams may notice failures informally ("this stage output tends to drift"). At Level 4, those observations become structured: test cases, regression checks, and evaluation criteria that can be applied whenever a template is revised. This is not about mathematical precision; it is about operational assurance. If a workflow is intended to never invent facts, the evaluation harness tests whether it invents facts under common stress scenarios. If a workflow is intended to stop when inputs are missing, the harness tests whether it actually stops. If a workflow is intended to avoid recommendation phrasing, the harness tests whether it produces recommendation language when prompted aggressively. Evaluation at Level 4 is therefore a governance tool: it turns "we think this works" into "we tested that it behaves within constraints."

Another Level 4 expansion is the management of internal IP. Workflow templates, prompt libraries, and QA checklists become proprietary assets. They encode the firm's way of working. That creates both opportunity and risk. The opportunity is consistent quality and faster onboarding. The risk is uncontrolled copying, uncontrolled modification, and unclear ownership. Level 4 governance therefore includes asset ownership, access control, and approved distribution. It also includes supervision frameworks: who is allowed to modify templates, how changes are reviewed, and how the organization ensures that a template remains aligned with professional standards and internal policies.

In short, Level 4 takes what Level 3 made possible (governed workflow runs) and turns it into a reusable system (governed workflow assets). The next chapter will show how to build, release, and maintain that system without drifting into automation bias or losing traceability.

---

**Artifact (Save This)**

**Level 3 exit criteria (ready to move to Level 4).**

1. The team uses defined workflow templates with stage boundaries and mandatory checkpoints.
2. Every workflow run produces a saved Workflow Run Bundle (logs, artifacts, reviewer notes, versions).
3. Stop conditions are enforced in practice (no silent progression when inputs are missing).
4. The team can identify and revise workflow templates based on observed failures (controlled change).

# Bibliography

[1] National Institute of Standards and Technology (NIST). *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*. NIST AI 100-1, U.S. Department of Commerce, 2023.

[2] International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC). *ISO/IEC 23894:2023 — Artificial intelligence — Guidance on risk management*. ISO/IEC Standard, 2023.

[3] Organisation for Economic Co-operation and Development (OECD). *Recommendation of the Council on Artificial Intelligence*. OECD Legal Instrument OECD/LEGAL/0449, 2019.

[4] European Parliament and Council of the European Union. *Regulation (EU) 2024/1689 of 13 June 2024 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act)*. Official Journal of the European Union, 2024.

[5] International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC). *ISO/IEC 42001:2023 — Information technology — Artificial intelligence — Management system*. ISO/IEC Standard, 2023.

[6] Board of Governors of the Federal Reserve System. *Supervisory Guidance on Model Risk Management*. Supervisory Letter SR 11-7, 2011.

[7] Committee of Sponsoring Organizations of the Treadway Commission (COSO). *Internal Control — Integrated Framework*. COSO Framework, 2013.

[8] International Organization for Standardization (ISO). *ISO 31000:2018 — Risk management — Guidelines*. ISO Standard, 2018.

# Chapter 4

# Innovators

**Abstract.** Chapter 4 introduces Level 4 (Innovators) in the AI Consulting governance-first maturity ladder. At this level, consulting firms move beyond one-off AI-assisted workflows and begin to design reusable internal assets such as prompt libraries, standardized templates, playbooks, and evaluation harnesses. The objective is not speed or automation, but disciplined innovation that reduces variance while preserving professional judgment. The chapter emphasizes that reuse fundamentally changes the risk profile: errors, hidden assumptions, and design flaws can scale across engagements if governance does not scale in parallel. Level 4 therefore requires explicit asset ownership, version control, documented assumptions, testing against known failure modes, and controlled release processes. Innovation is treated as internal consulting infrastructure, not as a productized substitute for expertise. All outputs remain drafts, all assets remain subject to human supervision, and no firm-wide automation is assumed. By the end of this chapter, readers will be able to design, evaluate, and govern reusable AI-enabled consulting assets in a way that is auditable, inspectable, and defensible, setting the foundation for organization-level adoption in Level 5.

---

**Artifact (Save This)**

**Scope disclaimer (required).** This chapter is an educational, governance-first guide for management consulting and corporate strategy. Level 4 focuses on *reusable AI-enabled assets*, not autonomous decision-making. All outputs remain drafts. Human judgment, supervision, and accountability are mandatory.

**Level alignment.** Level 4 introduces standardized templates, prompt libraries, evaluation harnesses, and controlled reuse. No firm-wide automation, no unsupervised deployment, and no client-facing autonomy.

## 4.1  Chapter overview: Level 4 in the maturity ladder

### 4.1.1  From workflows to assets

Level 3 and Level 4 can look deceptively similar from a distance because both involve multi-step work, repeated patterns, and a desire to reduce friction in delivery. The difference is not the surface form of the work; it is the *unit of value* the firm is trying to create. Level 3 is primarily about making consulting work *repeatable.* Level 4 is about making consulting work *reusable.* Repeatability means that a team can run a defined workflow again tomorrow with similar discipline: the same checkpoints, the same output schema, the same requirement to separate facts from assumptions, and the same human review gate before anything goes to a client. Reusability means something stronger and riskier: the firm is now building assets—templates, prompt libraries, evaluation harnesses, reusable storylines, diagnostic interview packs, memo shells, risk registers, and governance wrappers—that can travel across teams, practices, and engagements. In other words, Level 4 turns an internal workflow into a piece of internal IP.

This transition matters because the consulting firm is not merely accelerating one team's productivity; it is shaping the firm's *default behavior.* A workflow is typically owned by an engagement team and interpreted in context. An asset is designed to outlive the engagement, to be used by people who were not present when it was created, and to be applied under time pressure to problems that will never match the original conditions perfectly. This is exactly why assets create leverage in consulting and exactly why they create new governance obligations. The moment a firm says, "This is our standard AI-enabled market entry memo template," it is making a promise to itself: this asset will carry judgment-like structure into future work. If that structure is sloppy, biased, or internally inconsistent, the firm has not merely drafted a bad memo; it has encoded a bad habit.

A useful way to define the Level 3 to Level 4 transition is to say that Level 3 asks, "Can we run a supervised workflow with checkpoints and logs?" and Level 4 asks, "Can we *publish* the workflow as an internal asset without institutionalizing error?" Publishing is the key word. "Publishing" does not mean public release; it means internal release into the firm's operating environment: internal wikis, knowledge bases, onboarding kits, and practice toolkits. At that point, the firm is no longer governing only the *outputs*; it must also govern the *inputs to future outputs*, meaning the asset itself.

In consulting practice, this shift mirrors a familiar arc. Most firms begin with artisanal excellence: one partner, one team, a bespoke approach, and high-quality output driven by judgment and experience. Over time, the firm codifies what works: playbooks, checklists, interview guides, and standard storylines. Codification is how consulting scales. AI does not change that; it accelerates it. A good Level 4 program should therefore be described as "codification with stronger controls than we used before," not as "automation." The idea is not to let an AI system decide; the idea is to standardize the scaffolding around how humans decide, and to do so in a way that is inspectable and correctable.

Concretely, what changes when a firm moves from workflows to assets? First, the artifacts

become more formal. In Level 3, a workflow run bundle might include a structured prompt, the model output in JSON, human edits, reviewer notes, and the final deliverable. In Level 4, the firm must produce an *asset bundle* that includes not just one run but a description of intended use, explicit exclusions, embedded assumptions, a change log, and evaluation evidence. The asset bundle becomes the object that can be reviewed, approved, versioned, retired, or rolled back.

Second, the firm must adopt a clearer distinction between *engagement-specific context* and *asset-level defaults.* Engagement context is always incomplete and always contains client-specific nuance. Assets, by contrast, must be robust to missing data and must avoid implying certainty. This is why Level 4 assets should be designed to *surface unknowns* rather than to fill them. A high-quality consulting asset does not quietly invent a market size; it forces the consultant to either provide it or to log it as an open question and proceed with conditional reasoning. A high-quality asset does not produce a final recommendation; it produces a structured draft that makes the human reviewer's job easier: it clarifies what was provided, what was assumed, what remains open, and what would need to be verified before a client-facing statement can be made.

Third, roles and accountability become sharper. At Level 3, ownership can be reasonably localized: the engagement lead owns the workflow run, the reviewer signs off, and the artifacts are stored for reconstruction. At Level 4, an asset must have an owner independent of any single engagement: someone responsible for the asset's purpose, its scope boundaries, its performance under evaluation, its change history, and its retirement decision. This is not bureaucratic overhead; it is the minimum cost of reuse. Without it, assets sprawl, drift, and silently accumulate inconsistencies. In a consulting environment where teams already rely heavily on templates and precedent, silent drift is one of the most common ways a firm loses quality control.

Fourth, evaluation becomes a first-class deliverable. In Level 3, evaluation may be informal: the team notices the model sometimes invents facts, so they tighten the prompt. In Level 4, evaluation must be explicit because the firm is about to reuse the asset at scale. That evaluation does not have to be mathematically sophisticated; it has to be meaningful and documented. The purpose is not to prove truth; it is to measure predictable failure modes. For example, a storyline generator should be tested on cases with missing data, contradictory data, and ambiguous objectives. A diagnostic interview guide should be tested on multiple industries and operating models to confirm that it does not push a biased narrative. A market entry memo template should be tested against scenarios where the answer is "do not enter" and scenarios where the best option is "delay and investigate," so that the asset does not always produce optimistic momentum by default.

This is where the consulting mindset is helpful. Consultants already understand that a template is never neutral. A template guides attention: it privileges some questions and downplays others. A Level 4 AI asset is a template with a stronger rhetorical engine attached. It will produce fluent text that sounds plausible, and in doing so it will influence what teams perceive as "complete." This is why the chapter insists on governance-first design. A reusable asset must make it hard to produce confident claims without verification. It must make it easy to track what changed between versions. It must include explicit reminders that the output is a draft and that facts must be verified.

The most practical way to implement this in a consulting firm is to treat assets like internal products with a conservative release process. In many firms, internal templates are shared informally, copied, and modified. That culture is part of how knowledge spreads, but it is incompatible with high-stakes AI-enabled assets. Level 4 should introduce a controlled channel: assets are drafted, evaluated, reviewed, released with a version number, and then monitored. Teams can still iterate, but iteration becomes visible through versioning rather than hidden through informal copying. The goal is not to eliminate creativity; it is to ensure that the firm can answer, at any time, "Which asset version produced this draft, under what assumptions, and with what known failure modes?"

To make this concrete, consider a "one-page decision memo" asset for market entry. At Level 3, a team might run a workflow: gather inputs, prompt the model to draft, then review and edit. At Level 4, the firm creates a reusable memo asset: a structured prompt template that enforces the facts/assumptions/open questions schema; a standard memo skeleton (context, decision, options, risks, next steps); and an evaluation pack that tests the memo under three distinct patterns: data-rich entry, data-poor entry, and "do not enter." The asset includes explicit exclusions: it cannot be used to generate market sizing, it cannot cite sources, and it cannot present a final recommendation as fact. It includes an embedded rule: if the user provides no quantified inputs, the memo must present conditional ranges as assumptions and must request verification. This is a small design decision that becomes powerful at scale because it shifts the firm's default output away from false precision.

Or consider an operating model redesign workshop kit. At Level 3, a team might draft agendas and interview guides ad hoc. At Level 4, the firm creates an asset that standardizes workshop objectives, stakeholder mapping prompts, an agenda template with timeboxes, a decision-log structure, and a facilitation checklist. The asset is evaluated on different organizational contexts: matrixed global firms, product-led organizations, and regulated environments. The evaluation is not about "correct answers"; it is about whether the asset consistently surfaces the right unknowns and whether it avoids implying that a workshop output is a decision. The governance wrapper ensures that workshop notes are captured as decisions, actions, owners, and open items, and that unresolved topics are explicitly logged rather than quietly turned into "alignment."

This is the central promise of Level 4: not that AI will produce better judgment than consultants, but that it can help firms encode and distribute disciplined scaffolding that supports good consulting. In a firm that cares about quality, Level 4 is not optional experimentation; it is how the firm avoids reinventing governance on every engagement. But the transition must be treated as a structural change: the firm is now building internal machinery that will shape future client work. That machinery must be auditable, inspectable, and governable.

In practice, the firms that succeed at Level 4 are those that set expectations correctly. They do not call assets "solutions" or "engines." They call them "drafting scaffolds," "structured templates," "governance wrappers," and "evaluation-backed internal IP." They assume the model will sometimes be persuasive and wrong. They assume users will be rushed. They assume junior staff will over-trust fluent output. And they design assets that resist those pressures. This is why the ladder exists:

every increase in capability and scale requires a proportional increase in discipline. Level 4 is the first level where scale is the primary objective. That is why Level 4 is where many firms accidentally institutionalize their first serious failure. This chapter is designed to make that failure less likely by forcing a different mentality: from "Can we do it?" to "Can we publish it responsibly?"

### 4.1.2 Why Level 4 is dangerous without governance

Level 4 is the point at which a consulting firm can unknowingly turn a localized mistake into a repeatable institutional habit. The danger is not only that the model can be wrong. The danger is that the firm can become wrong *systematically*, with a professional tone, at a speed and scale that exceeds the firm's ability to notice. This is why Level 4 must be treated as a governance threshold. In Level 1 and Level 2, failures are often visible because they occur in individual drafts or reasoning structures. In Level 3, failures can be contained by stage gates: the workflow stops when an assumption register is incomplete or when human review has not happened. In Level 4, the failure mode is subtler: the asset itself becomes a distribution mechanism. People reuse it precisely because it looks legitimate, and because it carries the firm's implied endorsement. That implied endorsement is the risk multiplier.

The core problem is that reuse collapses context. A workflow run is attached to a specific engagement and a specific team. The team has situational awareness: they know which facts are firm, which are shaky, and which are politically sensitive. A reusable asset travels without that awareness. It must therefore be built to operate under uncertainty without pretending it has certainty. Without governance, assets will drift toward exactly the opposite behavior: they will become more polished, more confident, and more "complete" looking. That completeness is seductive in consulting because it resembles readiness. But readiness is not the same as correctness, and polished text can conceal both missing work and unverified claims.

A consulting firm should expect at least five categories of Level 4 danger, each arising from the interaction between reuse and human behavior.

The first danger is *institutionalized hallucination*. Many LLM failure modes are known: invented numbers, invented sources, invented causal links, and invented "industry norms." In one-off use, a sharp reviewer can catch these. In reuse, especially by junior teams under time pressure, invented details can survive because they appear inside a "standard template." A prompt that inadvertently invites the model to fill gaps ("make reasonable assumptions" without constraints) becomes a factory for plausible fiction. The firm will not see this as a model problem; it will see it as "how we write memos." That is the institutionalization risk. The firm's brand of professionalism becomes the carrier for unverified content.

The second danger is *assumption laundering*. Consulting work always involves assumptions. The ethical and professional requirement is to make assumptions explicit, owned, and testable. Without governance, reusable assets tend to hide assumptions in defaults: default market growth rates, default competitor sets, default "best practice" recommendations, default maturity model conclusions. These defaults become invisible because the asset produces them automatically. When

assumptions become invisible, they stop being debated. When they stop being debated, they stop being tested. The firm then begins to confuse convenience with truth. This is not only a technical problem; it is a professional risk because the firm can no longer demonstrate how it arrived at a conclusion. It can only demonstrate that it used a "standard tool."

The third danger is *false authority and decision laundering.* As assets become more integrated into delivery, there is a temptation to treat them as authoritative. A consultant may say, implicitly or explicitly, "This is our firm's standard approach," or "This is what the model produced," or "This is what the playbook recommends." In a client environment, those statements can suppress discussion. Inside the firm, they can suppress dissent. The danger is not that the asset is used; the danger is that the asset is used as an *excuse.* A firm that prides itself on judgment must resist the temptation to outsource responsibility to an internal tool. Level 4 is where that temptation becomes organizationally plausible because assets feel like firm property and therefore like firm endorsement.

The fourth danger is *silent drift.* Even if an asset starts out disciplined, it can degrade over time. Prompts are tweaked informally. Templates are copied and modified. People add "helpful" language. A junior consultant removes the open-questions section because it looks messy. A partner asks for a more confident tone. These changes are understandable in the moment, but over time they accumulate into a different asset than the one that was evaluated and approved. The firm then believes it is using a tested tool, but it is actually using a drifted variant. Drift is a governance failure because it breaks traceability: you cannot reconstruct what rules produced what output. In consulting, traceability is not merely a technical nice-to-have; it is how the firm protects itself when challenged. When a client asks why a claim was made, the firm must be able to show the source: the provided facts, the assumptions, and the review path. Drift destroys that path.

The fifth danger is *misuse outside intended scope.* Consulting assets are often designed for a particular kind of engagement: a market entry in consumer goods, a cost program in a manufacturing context, an operating model redesign in a matrixed organization. When the asset is reused elsewhere, users will often assume it is general-purpose. Without explicit scope boundaries, assets will be applied to contexts where their implicit assumptions are wrong. In a traditional template world, this is inconvenient. In an AI-enabled asset world, it becomes dangerous because the asset produces persuasive narrative that can overpower the user's doubt. The user is not merely copying a slide; they are being led through a structured argument that may not apply.

These dangers share a common theme: reuse creates *default behavior.* Governance is therefore not a compliance add-on; it is the mechanism by which the firm decides what its defaults should be. Level 4 governance can be summarized as answering three questions in a way the firm can defend.

First: *What is this asset allowed to do?* The asset must have explicit capabilities and explicit exclusions. If the asset is a storyline generator, it can structure narrative but it cannot invent data. If it is a diagnostic template, it can propose hypotheses but it cannot present them as findings. If it is a memo skeleton, it can produce a draft but it must label it as not verified and must surface open questions. These boundaries must be embedded not only in documentation but in the asset itself. Users do not reliably read documentation under pressure. The asset must enforce discipline

by design.

Second: *How do we know the asset behaves acceptably?* The answer cannot be "because it worked once." At Level 4, the firm needs an evaluation harness. The evaluation harness is not a claim of correctness; it is a claim of managed risk. It defines test cases and expected properties: no new facts, consistent separation of facts and assumptions, explicit open questions when data is missing, stable structure across different inputs, and predictable handling of ambiguity. The evaluation harness should also include adversarial tests: prompts that try to trick the asset into inventing citations, or into producing confident conclusions without inputs. The objective is to discover where the asset breaks and to document those breakpoints so users are warned and governance is tightened.

Third: *Who is responsible for the asset over time?* Without an owner, assets become abandoned. Abandoned assets are a liability because they are still used, but no one monitors them. Ownership must include a review cadence and a retirement policy. The firm must be willing to say, "This asset is deprecated; do not use it," and to provide a replacement path. This is standard product governance, but in consulting it often feels foreign because internal templates are treated casually. Level 4 requires the firm to treat AI-enabled assets as infrastructure, not as casual files.

There is also a more subtle governance issue that consulting leaders should anticipate: Level 4 can change the firm's culture of learning. In a healthy consulting culture, junior staff learn by struggling through problem structuring, building issue trees, conducting interviews, and synthesizing narrative. AI-enabled assets can accelerate learning if they are designed as scaffolds. But they can also suppress learning if they are treated as shortcuts. A Level 4 asset that always produces an "answer-shaped" memo can reduce the incentive to do the underlying thinking. Governance therefore includes not only risk controls but also pedagogical design: assets should be built to require the user to supply facts, to log assumptions, and to acknowledge uncertainty. This keeps the human in the loop not as a rubber stamp but as an active owner of the reasoning.

To keep the discussion grounded in consulting reality, consider how a small flaw becomes a firm-wide liability. Imagine a reusable "market overview" template that includes a section titled "Market size and growth." If the prompt does not explicitly forbid invention, the model will often fill that section with plausible numbers. One team catches it and deletes it. Another team does not notice. A third team uses the numbers as "directional" and repeats them in a client workshop. Now the firm has created a pattern: the template routinely manufactures market data. Even if the firm later corrects it, the damage is not just the incorrect numbers; it is the erosion of discipline. People begin to believe that "directional" invented data is acceptable if it sounds reasonable. That is the real institutional liability: the firm normalizes a lower standard of evidence.

The remedy is not to remove the section; the remedy is to govern it. The asset should require the user to provide the market size input explicitly, or else it should output a placeholder and list an open question: "Provide verified market size and growth (source and date)." It should include a reminder that the model does not research. It should place the burden back on the consultant to do the work or to mark it as pending. This is a small design choice that changes behavior at scale.

Another common example is the misuse of "best practices." Consulting language is full of them,

and clients often expect them. But "best practice" claims without evidence can be thinly disguised opinion. A reusable operating model playbook that encourages the model to produce best-practice statements will generate confident generalities that can crowd out engagement-specific nuance. Without governance, this becomes a style: the firm's assets sound smart but become generic. Over time, clients perceive the work as templated. The firm loses differentiation. This is a business risk as well as a governance risk. The control is to require assets to frame best practices as hypotheses with conditions: "In contexts where X holds, firms often do Y; verify whether X holds here." This retains usefulness while forcing explicitness.

It is also essential to recognize that Level 4 increases confidentiality and IP risks. Reusable assets often include examples, phrasing patterns, and distilled insights from past engagements. Without careful redaction and policy, assets can accidentally embed client-specific details or commercially sensitive approaches. Moreover, if the firm relies on external AI systems without appropriate controls, asset prompts themselves may leak internal methodologies. Governance at Level 4 therefore includes not only content discipline but also environment discipline: approved tools, approved data handling, and minimum-necessary inputs. The asset bundle should document these constraints so that reuse does not become an uncontrolled data channel.

Finally, Level 4 can tempt leaders into premature claims of standardization. A firm may want to tell itself that it now has a "platform" or a "capability." The chapter insists on resisting that narrative. Level 4 is innovation, not organizational transformation. The firm can and should build internal assets, but those assets remain drafts and scaffolds. They do not eliminate the need for engagement leadership, for partner judgment, for deep client listening, or for rigorous fact work. If leaders treat Level 4 assets as substitutes for expertise, they will create a brittle delivery model that fails precisely when the work becomes complex.

This is why the maturity ladder is governance-first. Each step increases leverage and therefore increases the cost of error. Level 4 is where leverage becomes a design goal. That is why it is dangerous without governance. Reuse will scale both value and error. A small flaw in a prompt can become a firm-wide habit. A hidden assumption can become a standardized conclusion. A persuasive but wrong narrative can become a "standard storyline." The only defensible response is to treat assets as governable objects: owned, evaluated, versioned, bounded, and auditable. That is the purpose of Level 4 in this book: to help consulting firms innovate without building a machine that amplifies their worst instincts.

In the chapters that follow, we will make this practical. We will define what an asset bundle must contain, what a minimal evaluation harness looks like, what change control means in day-to-day consulting life, and how to design assets so that they teach discipline rather than bypass it. If Level 3 was about proving that supervised workflows can exist, Level 4 is about proving that the firm can create reusable IP without surrendering responsibility. Level 5 will then extend this logic from assets to operating model: governance becomes not a chapter topic but a firm capability.

## 4.2 Mental model: what "innovation" means in consulting AI

Level 4 is the chapter where many consulting firms get the story wrong about what they are actually doing. They say they are "innovating," and what they often mean is that they are adopting a new tool. Or they say they are "scaling AI," and what they often mean is that they are moving from informal experimentation to something repeatable. Those are real changes, but they are not yet the real Level 4 shift. In this maturity ladder, Level 4 is not about adopting a model. It is about converting a set of supervised ways of working into reusable internal assets that shape how the firm behaves by default. The consultant is still responsible. The firm is still accountable. The AI system is still an unreliable generator of plausible text. What changes is that the firm begins to invest in internal infrastructure: curated templates, prompt libraries, playbooks, and evaluation harnesses that can be reused across many engagements without being reinvented each time.

This mental model matters because it determines what the firm optimizes for. If leadership thinks Level 4 is about "getting more output from the model," it will optimize for speed and volume. If leadership understands Level 4 as building internal infrastructure, it will optimize for inspectability, reuse boundaries, evaluation evidence, and long-term maintainability. In consulting, infrastructure is how quality scales. But infrastructure is also how mistakes scale. Level 4 therefore requires that innovation and governance be treated as inseparable. You cannot safely publish reusable AI-enabled assets without a control framework, because the asset becomes a silent teacher: it trains consultants on what "normal" looks like in draft work. If "normal" includes hidden assumptions, invented facts, and overconfident language, the firm will not only make errors; it will normalize them.

The most practical way to hold this in your head is to treat Level 4 as the point where the firm begins to curate a managed portfolio of "consulting accelerators" and to govern them like internal products. These accelerators are not client-facing products. They are not autonomous advisors. They are not decision engines. They are structured aids that reduce variance, standardize scaffolding, and make it easier for teams to produce disciplined drafts, consistent logs, and auditable workpapers. The more widely an accelerator is reused, the more it must behave like infrastructure: versioned, documented, tested, owned, and monitored. Level 4 is not "the model is smarter." Level 4 is "the firm is more disciplined about how it uses the model."

### 4.2.1 Useful abstraction

A useful abstraction for Level 4 is: *internal consulting infrastructure*. That phrase is deliberately unglamorous. It is meant to invoke the kinds of assets that consulting firms already rely on, but with the added reality that AI makes them more powerful and therefore more risky. Internal consulting infrastructure includes templates, libraries, playbooks, and evaluation harnesses. These assets do not replace consulting; they shape the structure within which consulting happens.

A template is a constraint on narrative. It dictates what sections exist, what must be stated, what must be left open, and what must be verified. A library is a curated set of reusable building blocks: prompt patterns for standard deliverables, tone variants, standard disclaimers, redaction

instructions, and structured output schemas. A playbook is a codified method: a diagnostic sequence, a set of interview questions tied to hypotheses, a workshop design, a synthesis pattern. An evaluation harness is the governance layer that tests whether the asset behaves acceptably under predictable conditions and failure modes.

This abstraction matters because it forces the right questions. When you treat an AI-enabled asset as infrastructure, you ask: Who maintains it? What is its scope? How is it versioned? How do we test it? What are the known failure modes? What happens if it drifts? Who can publish changes? How do we deprecate it? How do we train users? These questions are normal in software and product organizations. They are less normal in consulting, where templates often spread informally. Level 4 requires importing a small, conservative subset of product discipline into consulting knowledge management.

The abstraction also clarifies what "innovation" means in this book. Innovation is not novelty. Innovation is not publishing a new prompt every week. Innovation is building reusable assets that reduce error, reduce variance, and improve auditability. It is building scaffolds that make good work easier and bad work harder. It is turning tacit know-how into explicit structure without turning professional judgment into a canned answer.

Under this abstraction, the "AI" part becomes less mysterious. The model is simply one component in a system of assets. The most important design decisions are often not about model selection; they are about the constraints around usage. For example, a memo-drafting asset can be designed to always separate facts from assumptions, always list open questions, and always mark verification status as not verified. That design choice does more for governance than any change in model capability. Likewise, an operating model workshop kit can be designed to force decision logging and to separate "alignment statements" from actual decisions. Again, the model is not the source of discipline; the asset design is.

Finally, the infrastructure abstraction is useful because it encourages modularity. The firm should build small assets that can be combined rather than one giant "consulting engine." A storyline title generator can be a separate asset from a decision memo template. A redaction helper can be separate from a workshop agenda builder. A risk taxonomy checklist can be separate from a synthesis writer. Modularity makes inspection easier and misuse harder. It also makes evaluation practical: you can test each asset's specific failure modes rather than trying to validate a monolithic system.

In short, the useful abstraction is: Level 4 is the internal infrastructure layer of AI consulting. It is where the firm begins to create reusable assets, not to outsource thinking. The model is a drafting component; the asset is the disciplined wrapper around it.

### 4.2.2 Dangerous misconception

The dangerous misconception at Level 4 is that the firm has "productized judgment." This misconception appears in several forms, often dressed in optimistic language: "We have a standard approach now." "We have an AI playbook." "We've built a solution." "This is a proven asset." Each

phrase hides a category error: confusing a reusable scaffold for a validated answer.

In consulting, standardization is valuable. But standardization does not imply correctness. A standardized process can reliably produce wrong outputs if its assumptions are wrong or if its evaluation is weak. AI-enabled assets intensify this risk because they produce fluent text that looks complete. A firm that mistakes fluency for validation will begin to trust the asset as if it has been proven. In reality, many assets will only have been demonstrated once, in a narrow context, by the team that created them. That is not proof; it is anecdote.

A related misconception is: "The asset is proven because it worked once." This is the fastest way to institutionalize error. Many prompts and templates appear to work because they produce something plausible. They have not been tested against missing data, contradictory objectives, hostile inputs, or unusual industries. They have not been tested for drift. They have not been tested for the most common failure mode in consulting AI: inventing plausible but unverified claims. If an asset is reused widely before these failure modes are discovered, the firm is effectively running a large uncontrolled experiment on client work. Governance-first practice forbids that.

Another misconception is: "Standardization equals correctness." In a consulting firm, standardization often functions as a surrogate for quality. A clean template, a familiar storyline, and a crisp narrative create the appearance of rigor. With AI, this appearance becomes even stronger. The danger is that the firm becomes satisfied with the form of the deliverable and stops insisting on the substance: verified facts, explicit assumptions, and documented reasoning. The asset can become a factory for "answer-shaped" documents that are not supported by evidence. The client may not notice immediately, because the output looks like consulting. But the firm's risk profile changes: it is now producing consistent persuasive narratives that can be wrong in systematic ways.

There is also a governance misconception: "Because it's internal, it's safe." Internal assets can be more dangerous than external one-off drafts because they are reused and because they carry implied endorsement. If a junior consultant uses a random prompt they found, the firm can plausibly frame it as an individual error. If the same consultant uses an asset published in the firm's library, the firm has endorsed the structure. The firm is then accountable for the asset's behavior, even if no one intended that. Level 4 requires recognizing that internal publication creates institutional responsibility.

A final misconception is cultural: "Innovation means fewer constraints." In reality, Level 4 innovation requires more constraints, not fewer. Constraints are what make assets reusable without becoming liabilities. A constraint such as "do not invent facts" is not restrictive; it is liberating, because it makes reuse safer. A constraint such as "always output open questions when data is missing" is not bureaucratic; it is what prevents false precision from becoming standard. The firms that succeed at Level 4 innovate by designing constraints that preserve professional standards at scale.

The antidote to these misconceptions is to restate, repeatedly and plainly: a Level 4 asset is not an authority. It is a scaffold. It is a controlled template for producing drafts. It reduces variance in structure, not variance in truth. It is never "proven" in the sense of correctness; it is

only "evaluated" in the sense of known failure modes and documented boundaries. And it never removes responsibility from the human consultant.

### 4.2.3   Definition of a good Level 4 output

A good Level 4 output is one that can be reused across engagements *without smuggling in unverified content*, *without collapsing context*, and *without requiring hero-level vigilance* from reviewers to remain safe. It is reusable, documented, evaluated, owned, and supervised. It is easy to inspect and hard to misuse.

To make this definition operational in a consulting firm, it helps to describe Level 4 output as having three layers: the *asset layer*, the *run layer*, and the *human supervision layer*. The asset layer is the template or prompt library itself: its scope, exclusions, embedded assumptions, and output schema. The run layer is the instance of usage: the specific inputs provided in an engagement, the output produced, and the logged open questions and assumptions. The human supervision layer is the review and sign-off process: who reviewed, what was changed, what was verified, and what remains unverified.

A Level 4 asset is good if it produces run outputs that are predictably structured and governance-friendly. That means, at minimum, the asset must enforce the separation of facts provided, assumptions, open questions, and draft output. It must avoid creating new "facts." It must surface uncertainty rather than bury it. It must include explicit reminders of verification status. It must make it natural for a consultant to do the right thing: to request missing inputs, to mark assumptions, and to leave placeholders rather than inventing.

A Level 4 asset is also good if it is *evaluated* in a way that is meaningful for consulting practice. Consulting evaluation is not about numerical accuracy; it is about reliability of structure, stability of boundaries, and predictable behavior under missing information. A good asset comes with test cases that reflect real consulting conditions: incomplete inputs, conflicting stakeholder objectives, ambiguous scope, time pressure, and narrative temptations. The evaluation should document what the asset does well and where it fails. Failure is not disqualifying; undocumented failure is. The firm must know the asset's limits if it is going to reuse it responsibly.

Ownership is part of the definition of a good Level 4 output because without ownership the asset cannot remain good. The asset must have a named owner and a named reviewer. The owner is responsible for scope, evaluation, and updates. The reviewer is responsible for independent scrutiny. This is a minimal separation of duties. In consulting terms, think of the owner as the engagement manager for the asset and the reviewer as the partner or senior leader who signs off on the asset's publication. This is not about bureaucracy; it is about accountability for reusable IP.

A good Level 4 output is supervised not only at release but over time. It has a version number and a change log. It has a review cadence. It has a deprecation path. The firm should be able to answer, quickly: Which version is current? What changed? Why did it change? Who approved it? If an asset produced a problematic draft, can we trace that back to a version and decide whether to roll it back? These are infrastructure questions, and they become necessary the moment assets are

reused widely.

Finally, a good Level 4 output is hard to misuse. This does not mean users cannot misuse it; it means the asset is designed to resist common misuse patterns. It includes explicit exclusions so users are reminded not to use it for out-of-scope tasks. It includes prompts that refuse to invent facts. It requires users to provide inputs that matter. It outputs open questions when inputs are missing. It produces drafts that are clearly labeled as not verified. It does not present an "answer" when the available facts do not justify one. This is perhaps the most important property: a good Level 4 asset makes uncertainty visible.

The simplest way to test whether a Level 4 output is good is to perform a thought experiment. Imagine a new team member uses the asset under pressure, late at night, with incomplete data. Does the asset encourage them to invent, to overclaim, and to polish away uncertainty? Or does it encourage them to log assumptions, to surface open questions, and to preserve verification discipline? If the asset does the former, it is not a Level 4 asset; it is an institutional liability. If it does the latter, it is a candidate for safe reuse.

The artifact box below captures the minimum metadata required to treat an asset as governable infrastructure rather than as a clever prompt. This metadata is the firm's defense against drift, misuse, and false authority. It is also the firm's mechanism for learning: by forcing explicit scope, evaluation, and ownership, the firm can iterate responsibly. Without these fields, the firm cannot reconstruct what happened when an asset produces a problematic output. That inability is the hallmark of unmanaged risk.

---

**Artifact (Save This)**

**Non-negotiable rule. Facts are not assumptions.** Facts must be provided or verified; assumptions must be stated, owned, and testable.

**Required metadata for Level 4 assets (minimum).**

a. asset_name and purpose
b. scope and explicit exclusions
c. assumptions embedded
d. evaluation criteria and test cases
e. owner and reviewer
f. version and change log
g. verification_status: `"Not verified"`

---

A final note on tone: Level 4 outputs should avoid the temptation to sound like finished consulting deliverables. Their job is to make drafting disciplined and reviewable, not to mimic certainty. A firm that designs Level 4 assets to maximize polish will maximize risk. A firm that designs Level 4 assets to maximize inspectability will maximize defensibility. In consulting, defensibility is not a legalistic concept; it is a professional standard: can you explain what you did, what you assumed, what you verified, and why your recommendation is justified? Level 4 assets should make that

standard easier to meet, not easier to dodge.

This is why the mental model is so important. Innovation at Level 4 is not about producing more content. It is about building internal infrastructure that shapes the firm's default behavior toward disciplined drafting, explicit uncertainty, and auditable reuse. The model is not the authority. The asset is not the authority. The consultant is the authority, and the firm's governance is the mechanism that keeps that authority honest.

## 4.3  What Level 4 can do, and what it cannot do

Level 4 is where consulting firms begin to feel a new kind of confidence: "We are not just experimenting anymore; we are building something reusable." That confidence is partly justified. When done well, Level 4 can reduce wasted effort, improve consistency in drafting, and make governance practices easier to apply across many teams. But Level 4 is also where firms are most likely to confuse *reusability* with *authority*. A reusable asset carries the firm's implied endorsement. That endorsement is exactly why Level 4 produces leverage, and it is exactly why Level 4 must be explicit about boundaries.

A practical way to read this section is to treat it as a contract between leadership and the firm. The "CAN" statements define what Level 4 is for: the disciplined creation of reusable internal infrastructure that supports consulting work. The "CAN'T" statements define what Level 4 must never be allowed to become: a mechanism for outsourcing responsibility, laundering judgment, or quietly bypassing professional controls. In a governance-first approach, the firm does not earn the right to expand capability by building clever tools; it earns that right by building controls that scale with reuse.

This is why Level 4 must be framed carefully in internal communication. If the firm tells teams, "We are rolling out AI tools to speed up delivery," it will get speed—and also hidden assumptions, persuasive errors, and uncontrolled drift. If the firm tells teams, "We are publishing reusable assets to standardize structure and strengthen governance," it will get slower adoption—and safer scale. Level 4 is a choice about what the firm optimizes for. Consulting firms that treat Level 4 as a productivity initiative tend to create liabilities. Firms that treat Level 4 as an internal infrastructure initiative tend to create durable capability.

The distinction between CAN and CAN'T is not philosophical. It is enforceable in asset design. A well-governed asset can make some behaviors easy (structured drafting, explicit open questions, standardized storylines) and make other behaviors difficult (inventing facts, hiding assumptions, issuing recommendations without review). Level 4 success is not just a matter of policy; it is a matter of designing assets that embed the policy in usage.

### 4.3.1  CAN: innovation with discipline

At Level 4, consulting firms can innovate in ways that are both meaningful and governable, but only if they define innovation as "reusable scaffolding under control" rather than "new outputs at

scale." The capabilities below are the right targets because they create value through structure, not through false claims of correctness. They also align naturally with consulting work, where much of the craft is about narrative architecture, decision hygiene, and disciplined synthesis.

- **Create reusable prompt libraries for common consulting tasks.**

A consulting firm already has recurring deliverables: decision memos, weekly updates, stakeholder emails, interview guides, workshop agendas, diagnostic summaries, and slide storylines. A Level 4 prompt library formalizes the best prompt patterns the firm has learned so far and publishes them as internal assets with explicit governance wrappers. The core value is not that the library produces "better content." The core value is that it reduces variance in structure and discipline, especially under time pressure.

A properly designed prompt library does several things consistently. It enforces the separation of facts provided, assumptions, and open questions. It forbids the invention of numbers, benchmarks, citations, or "industry norms" unless the user supplies them. It standardizes disclaimers and verification status. It includes explicit scope boundaries: what the prompt is intended for and what it must not be used for. It also makes reuse auditable by including an asset name and version field so teams can log which asset produced which draft.

The most valuable consulting prompt libraries are not "general writing prompts." They are role-specific and deliverable-specific. They encode the firm's standards for what a good draft looks like: not merely polished, but inspectable and reviewable. A firm that publishes a prompt library is therefore publishing a standard of practice. This is why Level 4 requires ownership, evaluation, and change control around the library. Without those controls, the library becomes a chaotic collection of clever prompts that drift, contradict each other, and invite misuse.

- **Standardize slide storylines, memo structures, and workshop designs.**

Consulting quality is often driven by structure. A well-structured memo makes it harder to hide uncertainty. A well-structured storyline clarifies what is known, what is assumed, and what must be decided. A well-designed workshop forces decision logging rather than producing vague "alignment." These are exactly the domains where Level 4 assets can create disproportionate value because they shape the firm's default behavior.

Standardizing storylines does not mean standardizing conclusions. It means standardizing the shape of argument: a clear statement of the decision, a small number of options, explicit tradeoffs, risks, and next steps. Standardizing memo structures means requiring explicit sections for assumptions and open questions, not burying those elements in footnotes. Standardizing workshop designs means embedding governance into facilitation: agenda templates that include decision points, pre-read expectations, and a structured decision record produced at the end.

The key governance requirement is that these standardized structures must be designed to *resist false completeness*. Storylines should include placeholders where data is missing. Memo templates should require sources for quantitative claims. Workshop kits should explicitly distinguish between hypotheses, observations, decisions, and actions. Without these constraints, standardization

becomes a way to produce persuasive artifacts quickly, which is precisely how consulting firms can drift into "templated confidence" rather than defensible work.

A good Level 4 storyline asset, for example, should be "titles only" by default, with explicit tags where evidence is required. A good memo asset should include an explicit section titled "What we do not yet know" and require that section to be populated. A good workshop asset should require a decision log format that captures decisions, owners, and dates. These are small design choices that create governance at scale.

- **Build evaluation harnesses to test consistency and failure modes.**

Evaluation harnesses are one of the most important innovations at Level 4 because they operationalize humility. They are the firm's mechanism for saying: "We do not trust this asset because it looks good. We trust it only to the extent that we have tested its behavior and documented its limits." This is a fundamentally governance-first stance.

A practical evaluation harness for consulting assets does not need to be complex. It needs to be disciplined. It should include representative test cases that mirror real engagement conditions: incomplete inputs, contradictory stakeholder objectives, ambiguous scope, and adversarial prompts that try to induce invention. The harness should test for properties that matter in consulting governance: no new facts, stable separation of facts and assumptions, explicit open questions when data is missing, and refusal to provide final recommendations without stated assumptions and human review.

Evaluation should also test for drift-prone behaviors. Does the asset become more confident when inputs are thin? Does it quietly introduce "best practices" that are not tied to provided facts? Does it produce numbers when asked to "estimate"? Does it cite sources when none were provided? These are not academic concerns; they are the predictable ways in which AI outputs can become persuasive liabilities.

A firm that builds evaluation harnesses can safely iterate. Without them, iteration becomes uncontrolled change. The harness also creates a common language for governance discussions: instead of arguing about whether an asset "feels good," teams can review documented failure modes and decide whether the asset is safe for internal release.

- **Reduce variance in junior team outputs without removing review.**

One of the most legitimate motivations for Level 4 is variance reduction. Consulting firms are apprenticeship organizations. Quality depends on training, supervision, and repetition. Junior staff often struggle with structure: how to write a decision memo that is not a narrative dump; how to build a storyline that has a "so what"; how to run a workshop that produces decisions rather than discussion. Level 4 assets can provide scaffolding that helps juniors produce better drafts faster.

But governance-first practice requires a clear line: reducing variance is not the same as removing supervision. Level 4 assets should be designed as training wheels, not autopilots. They should make it easier for juniors to do the right thing, but they should still require human review, and they should still force explicitness about assumptions and open questions.

In fact, the best Level 4 assets make review more effective. They create consistent sections so reviewers can focus quickly on the highest-risk areas: What facts were provided? What assumptions are doing the work? What open questions remain? Where is verification needed? Review becomes less about rewriting and more about challenging the reasoning and evidence. This is exactly how Level 4 should improve consulting delivery: by shifting effort from formatting to judgment.

The discipline requirement is that juniors must not be able to use assets as a shield. If a junior says, "This is what the tool generated," that is a governance failure. The correct posture is: "This is a draft scaffolded by the asset; here are the facts I provided, here are the assumptions it forced me to state, and here are the open questions we must verify." Level 4 assets should be designed to support that posture, not to undermine it.

### 4.3.2   CAN'T: boundaries you must enforce

The boundaries at Level 4 are not optional guidelines. They are the difference between internal innovation and institutionalized irresponsibility. In consulting firms, the strongest temptation is to treat a reusable asset as a shortcut around senior judgment. The firm must actively resist that temptation by enforcing boundaries in three places: in policy, in asset design, and in engagement practice.

- **Replace partner judgment or engagement-specific reasoning.**

Consulting judgment is contextual. It depends on the client's strategy, constraints, politics, culture, risk appetite, and capabilities. No reusable asset can capture this context reliably, and no AI-generated output should be treated as a substitute for engagement leadership. Level 4 assets can structure drafts and surface questions, but they cannot decide what the client should do, what tradeoffs matter most, or what risks are acceptable.

This boundary must be reinforced culturally. Partners and senior leaders must treat assets as aids, not as authorities. If leadership begins to treat asset outputs as "standard answers," juniors will follow. The result is decision laundering: the firm quietly replaces judgment with templated narrative. This is precisely what governance-first practice forbids.

- **Auto-generate client recommendations without review.**

This boundary is both obvious and frequently violated in practice, often unintentionally. A storyline generator can produce slides that look like recommendations. A memo template can produce a conclusion section that reads like a recommendation. A workshop kit can produce "next steps" that sound like commitments. If these outputs are sent to clients without review, the firm has effectively allowed an internal asset to speak externally.

Level 4 therefore requires a strict human review gate for any client-facing communication influenced by AI assets. The firm should treat "review" as a real act: checking facts, challenging assumptions, confirming that the narrative is defensible, and ensuring that uncertainty is not being polished away. The firm should also require documentation of review: reviewer initials, date, and

a brief note on what was verified or changed. The objective is not to slow work; it is to preserve accountability.

- **Remove accountability through "standard tool" arguments.**

  One of the most dangerous governance failures is when a firm begins to treat internal assets as a way to diffuse responsibility. If a deliverable is challenged, someone says, "That's the standard template." If an assumption is questioned, someone says, "That's what the tool uses." These arguments are unacceptable in a governance-first model because they invert accountability: instead of humans owning the work, the tool becomes the scapegoat or the authority.

  To prevent this, Level 4 assets must be framed as *assistive scaffolds* with explicit ownership and explicit limits. Users must be trained to treat them as drafts. Reviewers must be trained to challenge them. And leadership must be trained to reject "standard tool" excuses. In a consulting firm, accountability is not transferable. If an asset is standard, that increases accountability; it does not reduce it.

- **Bypass confidentiality, conflict, or independence requirements.**

  Level 4 increases the risk of confidentiality and compliance failures because assets encourage reuse and because reuse often involves copying and sharing. Teams may paste sensitive client information into prompts to "make the asset work better." They may embed examples from past engagements into asset templates. They may store outputs in locations that are not approved for client-sensitive material. They may reuse a diagnostic playbook in a context where independence or conflict checks should have limited the work.

  Governance-first practice requires treating these risks as structural. The controls must not rely solely on individual judgment. Assets must include redaction instructions and minimum-necessary input guidance. Asset documentation must specify approved environments and prohibited data types. Where the firm has formal independence or conflict policies, asset usage must not become a workaround. In a consulting firm, the reputational damage from confidentiality failures is often greater than the direct operational damage. Level 4 cannot be allowed to become a new leakage channel.

  The central message of this section is that Level 4 is not a license to move faster. It is a commitment to move *more consistently* and *more defensibly*. What Level 4 can do is build internal infrastructure that helps teams draft better, structure thinking, and reduce variance. What Level 4 cannot do is replace judgment, bypass review, diffuse accountability, or weaken professional controls. If the firm enforces these boundaries, Level 4 becomes a platform for safer scale. If the firm fails to enforce them, Level 4 becomes the point where the firm institutionalizes persuasive error and calls it innovation.

## 4.4 The Level 4 governance wrapper (minimum controls)

---

**Risk & Control Notes**

**Capability.** Reusable AI-enabled consulting assets

**Primary risks.** Institutionalized error; silent drift; misuse outside intended scope; false authority; IP leakage

**Minimum controls.** Asset registry; version control; evaluation harnesses; usage boundaries; named owners; review cadence

---

Level 4 is the first point in the maturity ladder where governance can no longer be treated as an overlay on top of good intentions. In Levels 1 through 3, governance can often be enforced through discipline at the level of individual runs: structured prompts, assumption registers, human review gates, and saved artifacts. At Level 4, that approach breaks down. The object being governed is no longer a single output or workflow run; it is a reusable asset that will shape behavior across engagements, teams, and time. This requires a different kind of control framework: one that treats assets as governable objects in their own right.

The phrase "governance wrapper" is deliberate. It signals that governance is not a separate policy document or a training slide deck. It is the structure that surrounds the asset and constrains how it is created, evaluated, published, used, modified, and retired. Without this wrapper, Level 4 innovation collapses into informal template sharing, which is exactly how most consulting firms already operate. AI simply amplifies the consequences of that informality. With AI-enabled assets, a small design flaw can be reproduced hundreds of times with persuasive confidence. The governance wrapper exists to prevent that amplification from becoming institutionalized error.

At a high level, the Level 4 governance wrapper must accomplish five things simultaneously. First, it must make assets visible: the firm must know what assets exist, who owns them, and where they are used. Second, it must make asset behavior inspectable: the firm must be able to understand what an asset is designed to do, what assumptions it embeds, and where it is likely to fail. Third, it must control change: updates to assets must be intentional, documented, and reversible. Fourth, it must enforce boundaries: assets must not be silently repurposed outside their intended scope. Fifth, it must preserve accountability: no asset may become a proxy for judgment or a shield against responsibility.

Each of these objectives maps directly to the risks identified in the capability–risk–control triad. Institutionalized error arises when assets are reused without evaluation or oversight. Silent drift arises when assets change informally over time. Misuse outside scope arises when assets are not bounded. False authority arises when assets are treated as validated answers rather than scaffolds. IP leakage arises when reusable prompts and templates are not governed as sensitive internal material. The minimum controls listed in the risk box are not aspirational best practices; they are the minimum conditions under which Level 4 reuse is defensible in a professional consulting firm.

The rest of this section explains how those controls work together as a coherent wrapper, and

why each one is necessary.

## Why a wrapper, not a checklist

It is tempting to treat governance as a checklist: register the asset, assign an owner, test it once, and move on. That temptation should be resisted. Governance at Level 4 is not about one-time compliance; it is about lifecycle control. A wrapper is more than a checklist because it defines how the asset lives inside the firm over time. It answers questions that only arise months after an asset is created: Who updated this? Why was it changed? Which engagements used the previous version? Is this asset still appropriate for current practice priorities? Should it be retired?

In consulting firms, templates and playbooks often escape lifecycle control because they are treated as static knowledge artifacts. AI-enabled assets cannot be treated that way because their behavior depends not only on static text but on interaction with a probabilistic system. The same prompt can behave differently under slightly different conditions. That variability makes lifecycle governance essential. The wrapper provides a stable frame around an unstable core.

## Asset registry: making reuse visible

The first control in the Level 4 wrapper is an asset registry. Without a registry, the firm cannot even answer the most basic governance question: "What reusable AI-enabled assets do we have?" Informal sharing through chat channels, shared drives, or personal notebooks is incompatible with Level 4. The registry does not need to be complex. It does need to be authoritative.

At minimum, the asset registry should record the asset name, a short description of its purpose, its intended use cases, explicit exclusions, the current version, the owner, and the last review date. This turns assets from invisible tools into visible infrastructure. Visibility is the precondition for accountability. If an asset causes a problem and no one knows it exists, governance has already failed.

The registry also plays a cultural role. It signals that assets are firm property, not personal hacks. That shift matters because it changes how people treat the asset. When something is in the registry, it carries implied endorsement and therefore implied responsibility. Teams become more careful about publishing assets, and leaders become more careful about approving them. This friction is healthy. Level 4 is not about maximizing the number of assets; it is about maximizing the quality of assets that survive governance.

## Version control: governing change, not freezing it

Once an asset is visible, the next challenge is change. Consulting work evolves. Practices refine their thinking. New failure modes are discovered. Prompts that worked well six months ago may need adjustment. The governance wrapper must therefore allow change without allowing drift.

Version control is the mechanism that makes this possible. Every Level 4 asset must have a version identifier and a change log that explains what changed and why. This does not require

sophisticated tooling. It requires discipline. A simple versioning scheme (v1.0, v1.1, v2.0) and a short rationale for each change is sufficient. What matters is that changes are intentional and traceable.

Version control serves three critical functions. First, it allows the firm to reconstruct history. If a problematic output appears, the firm can determine which asset version produced it and what assumptions were embedded at the time. Second, it enables rollback. If a change introduces a new failure mode, the firm can revert to a previous version while investigating. Third, it enforces review discipline. A version change should trigger review by the named reviewer, not be merged silently.

Without version control, silent drift becomes inevitable. Small "improvements" accumulate. Defaults change. Assumptions shift. Over time, the asset becomes something different from what was evaluated and approved. In a consulting firm, this drift is particularly dangerous because it is rarely malicious. It is the product of good intentions under time pressure. Governance must be designed for realistic human behavior, not ideal behavior. Version control is how the firm protects itself against its own helpfulness.

## Evaluation harnesses: documenting known failure modes

Evaluation is the most misunderstood control in Level 4. Many firms hear "evaluation" and think of benchmarking or accuracy testing. That is not the objective here. The objective is to document how the asset behaves under conditions that matter for consulting practice.

A Level 4 evaluation harness should be designed around failure modes, not success stories. It should ask: What happens when inputs are missing? What happens when assumptions conflict? What happens when the user asks for something out of scope? Does the asset invent facts? Does it become overly confident? Does it blur the line between draft and recommendation? These are the questions that determine whether an asset is safe to reuse.

The evaluation harness should be saved as part of the asset bundle. This is not busywork. It is the firm's evidence that it exercised professional care before publishing a reusable asset. When leadership later asks, "Why did we think this was safe to use?" the evaluation results are the answer. They show that the firm did not assume correctness; it tested for predictable failure.

Importantly, evaluation does not certify correctness. It certifies awareness. A good evaluation document does not say, "This asset is correct." It says, "Here is what this asset does reliably, here is where it fails, and here are the conditions under which it must not be used." That distinction is fundamental. Consulting firms trade on judgment, not on guarantees. Evaluation harnesses support judgment by making risk visible.

## Usage boundaries: preventing scope creep

Even a well-designed asset can become dangerous if it is used outside its intended context. Usage boundaries are therefore a core part of the governance wrapper. Every Level 4 asset must have explicit scope and explicit exclusions. These boundaries must be documented and, where possible, enforced by the asset itself.

For example, a market entry memo template may be appropriate for early-stage strategic discussion but not for final investment committee decisions. A cost diagnostic prompt may be appropriate for hypothesis generation but not for setting savings targets. An operating model workshop kit may be appropriate for internal alignment but not for regulatory-facing documentation. These distinctions must be made explicit.

Usage boundaries serve two purposes. First, they protect the firm from misuse. Second, they protect the asset from being blamed for things it was never designed to do. When an asset is challenged, the firm should be able to say, "This asset was not intended for that use, and here is the documented boundary." Without such boundaries, every asset becomes a general-purpose tool by default, which is precisely how false authority emerges.

Boundaries should be treated as part of the asset's interface, not as fine print. They should appear in the asset description, in the prompt text, and in the output disclaimers. Users under pressure do not read documentation carefully. Governance must therefore be embedded in the interaction, not merely appended to it.

## Named owners and reviewers: anchoring accountability

Ownership is the most human control in the Level 4 wrapper, and also the most important. An asset without an owner is an orphan. Orphan assets are still used, but no one feels responsible for their behavior. This is how institutional risk accumulates silently.

Every Level 4 asset must have a named owner and a named reviewer. The owner is responsible for the asset's purpose, scope, evaluation, and ongoing maintenance. The reviewer provides independent scrutiny at publication and at major version changes. This is a minimal separation of duties that mirrors familiar consulting roles: engagement manager and partner reviewer.

Ownership does not mean that the owner personally approves every use of the asset. It means that the owner is accountable for the asset's design and governance. If the asset drifts, the owner is responsible for correcting it. If the asset proves unsafe, the owner is responsible for deprecating it. This clarity is essential. Without it, assets become everyone's problem and therefore no one's problem.

The reviewer role is equally important. Review provides a check against local bias. Asset creators are often too close to their work to see its weaknesses. A reviewer can challenge assumptions, question scope, and insist on clearer boundaries. In consulting firms, where senior review is already part of the culture, this role should feel familiar. Level 4 simply formalizes it for internal assets.

## Review cadence: governing assets over time

Governance does not end at publication. Assets must be reviewed periodically to ensure they remain appropriate. Changes in firm strategy, industry focus, regulatory environment, or AI tooling can all affect asset suitability. A review cadence is the mechanism by which the firm acknowledges that assets age.

The review cadence does not need to be frequent. Annual or semi-annual review is often sufficient. What matters is that review is scheduled and recorded. During review, the owner and reviewer should ask: Is this asset still used? Is it used in the contexts we intended? Have new failure modes emerged? Should it be updated, restricted, or retired?

Retirement is an underappreciated governance act. An asset that is no longer appropriate should be explicitly deprecated. Quiet abandonment is dangerous because old assets often continue to circulate. A clear deprecation notice, recorded in the registry, is part of responsible governance.

### 4.4.1 The Asset Bundle (what you must save)

All of the controls described above converge in a single concrete artifact: the Asset Bundle. The Asset Bundle is the minimum deliverable that turns an AI-enabled template or prompt into a governable asset. Without an Asset Bundle, the firm cannot demonstrate that it exercised care in publishing reusable infrastructure.

The Asset Bundle is not a bureaucratic artifact. It is the firm's memory. It captures the intent, the design, the evaluation, and the accountability around the asset. When questions arise later, the Asset Bundle is what allows the firm to answer them.

---

**Artifact (Save This)**

**Minimum deliverable standard (Asset Bundle).**

1. **Asset description and intended use cases.** A concise statement of what the asset is for, who should use it, and in what contexts. This anchors scope and prevents accidental generalization.
2. **Prompt templates or workflow definitions.** The actual reusable material, written clearly and versioned. This is the operational core of the asset.
3. **Evaluation results (test cases + failure modes).** Documentation of how the asset was tested, what predictable failures were observed, and what boundaries were set as a result.
4. **Assumption register embedded in the asset.** A clear list of assumptions the asset relies on, either by default or by design. This prevents assumption laundering.
5. **Version history and change rationale.** A record of how the asset has evolved and why changes were made. This protects against silent drift.
6. **Named owner and reviewer sign-off.** Explicit accountability for the asset's design and publication.

---

The Asset Bundle should be stored alongside the asset in the firm's knowledge system. It should be easy to find and easy to inspect. Teams should be trained to expect it. Leadership should insist on it. Over time, the existence of Asset Bundles becomes a signal of maturity: a firm that can show its Asset Bundles can demonstrate that its AI innovation is governed, not improvised.

**Why these are minimum controls**

It is worth emphasizing that the controls described in this section are minimums, not best-in-class aspirations. They are the least a consulting firm must do to justify reusable AI-enabled assets. A firm may choose to add more controls: usage analytics, automated compliance checks, or formal risk ratings. Those additions may be appropriate in some contexts. But without the basics—registry, versioning, evaluation, boundaries, ownership, and review cadence—additional sophistication is cosmetic.

Level 4 is often where firms are tempted to over-engineer. They build complex platforms but neglect basic discipline. This chapter takes the opposite stance. It argues that a small number of well-chosen controls, consistently applied, is sufficient to manage most of the risk. The challenge is not technical difficulty; it is organizational follow-through.

In consulting, credibility is cumulative. Clients trust firms that can explain how they work, not just what they deliver. The Level 4 governance wrapper is part of that explanation. It allows the firm to say, honestly: "We innovate, but we do so with control. We reuse assets, but we do not outsource judgment. We benefit from AI, but we do not hide behind it." That posture is not only safer; it is also strategically sound. In a world where clients are increasingly skeptical of opaque AI claims, governance-first innovation becomes a differentiator.

Level 4 governance is therefore not an internal tax on creativity. It is the mechanism that allows creativity to scale without undermining professional standards. Firms that embrace this will find that their assets become better over time, not just more numerous. Firms that ignore it will eventually find themselves explaining why a "standard tool" produced something that no one can defend.

## 4.5 Core Level 4 asset patterns (A–D)

By the time a consulting firm reaches Level 4, it should resist the temptation to invent dozens of bespoke AI-enabled tools. Proliferation is not innovation; it is unmanaged complexity. The purpose of Level 4 is not to maximize the number of assets, but to identify a small set of reusable patterns that can cover a large share of consulting work while remaining governable. Over time, most successful Level 4 programs converge on a limited number of asset archetypes. These archetypes are not defined by technology; they are defined by how consulting work actually happens.

The four patterns described in this section represent a practical and conservative starting set. Together, they address the core needs of consulting delivery at scale: disciplined drafting, structured synthesis, controlled evaluation, and safe skill development. They are deliberately framed as patterns rather than tools because the implementation details will vary by firm, practice, and engagement type. What should not vary is the governance posture embedded in each pattern.

Each pattern is designed to produce leverage without creating false authority. Each one can be evaluated independently. Each one can be owned, versioned, and retired. And critically, each one reinforces the central principle of this book: AI assists consulting judgment; it does not replace it.

### 4.5.1 Pattern A: Prompt libraries for recurring consulting tasks

Prompt libraries are often the first Level 4 asset firms attempt to build, and for good reason. Consulting work is highly repetitive at the level of form, even when it is bespoke at the level of content. Emails, memos, meeting summaries, decision records, diagnostic outlines, and workplans appear in almost every engagement. A prompt library captures the best known ways to scaffold these drafts and makes them reusable across teams.

The key to a governance-first prompt library is that it must be *deliverable-specific*, not generic. A library of "writing prompts" is not a Level 4 asset; it is a productivity hack. A Level 4 prompt library is organized around consulting artifacts that already have professional standards attached to them. For example: a one-page decision memo, a market scan summary, a cost diagnostic outline, a PMO weekly update, or a workshop follow-up email. Each prompt encodes not just tone, but structure and discipline.

A well-designed prompt library does four things consistently. First, it constrains scope. Each prompt states explicitly what it is intended to produce and what it must not be used for. A market scan memo prompt, for instance, may explicitly forbid market sizing, benchmarking, or citation unless the user provides those inputs. Second, it enforces separation of facts, assumptions, and open questions. This is not optional decoration; it is the core governance mechanism. Third, it standardizes output structure so that reviewers know where to look for risk. Fourth, it embeds verification status and disclaimers so that no draft can be mistaken for a final deliverable.

Consider a prompt library entry for a cost diagnostic outline. At Level 3, a team might write an ad hoc prompt to generate a list of cost levers. At Level 4, the firm publishes a prompt asset that does something more disciplined. It asks the user to provide the scope (business unit, geography, time horizon), known constraints (labor agreements, regulatory limits), and available data. It then produces an outline with sections for cost categories, hypotheses, required data, risks, and open questions. Crucially, it does not estimate savings or rank initiatives unless the user supplies explicit inputs. The output is useful precisely because it refuses to overreach.

Prompt libraries also benefit from being modular. Instead of one giant prompt that does everything, the library should contain smaller, composable prompts. One prompt drafts a memo structure. Another rewrites tone variants. Another extracts open questions. Another produces a decision log. Modularity improves governance because each component can be evaluated and versioned independently. It also improves learning: users understand what each prompt is responsible for, rather than treating the library as a black box.

From a governance perspective, prompt libraries should be treated as internal IP. They should live in a controlled repository, be versioned, and have named owners. Changes to widely used prompts should trigger evaluation and review. The firm should resist the urge to allow everyone to "improve" prompts informally. That behavior feels collaborative, but it is how silent drift enters the system. Innovation at Level 4 is not crowdsourcing prompts; it is curating them.

Finally, prompt libraries should be accompanied by usage guidance that is brief but explicit. Consultants under pressure will not read long manuals. The guidance should answer three questions:

When should I use this prompt? When should I not use it? What must I verify before anything leaves the team? If the prompt library cannot answer those questions clearly, it is not ready for Level 4 reuse.

### 4.5.2 Pattern B: Standardized storyline generators

Storylines are one of the most powerful—and dangerous—assets in consulting. A good storyline shapes how a problem is understood. It directs attention, frames tradeoffs, and implicitly signals what matters. This is why storyline generators are a natural Level 4 asset and why they require particularly strong governance.

The safest and most effective form of storyline generator at Level 4 is *titles-only*. That constraint is not arbitrary. Titles define the logic of an argument without pretending to provide evidence. They make gaps visible instead of filling them. A titles-only storyline generator helps consultants structure thinking while preserving the obligation to supply facts, analysis, and judgment separately.

A Level 4 storyline generator should therefore be designed to produce slide titles and section headers, not body text and certainly not numbers. It should include explicit placeholders where data is required and explicit markers where assumptions are being made. For example, a market entry storyline generator might produce titles such as "Market attractiveness depends on verified demand growth (TBD)" or "Economics sensitive to assumed margin range (to be validated)." These titles guide discussion without laundering uncertainty.

The governance value of standardized storylines lies in consistency. When all teams use a similar storyline structure for a given class of problem, reviewers can focus on substance rather than format. They know where to find the decision, the options, the risks, and the open questions. This reduces review time and improves quality. But consistency must not become rigidity. A good storyline generator should allow optional branches and explicitly state when a section may not apply.

Storyline generators also need explicit scope boundaries. A storyline designed for internal strategy discussion should not be reused for board materials without modification and review. A storyline designed for hypothesis generation should not be mistaken for a recommendation framework. These boundaries should be encoded in the asset description and reinforced in the output disclaimers.

Evaluation of storyline generators should focus on narrative failure modes. Does the generator always push toward a positive conclusion? Does it default to "best practices" language? Does it collapse multiple uncertainties into a single confident title? Does it obscure dissenting views? These are subtle risks, but they are exactly the risks that can reshape a firm's culture if left unchecked. A storyline generator that always produces momentum-oriented narratives can quietly bias the firm toward action even when caution is warranted.

Used well, standardized storyline generators can elevate thinking across the firm. Used poorly, they can homogenize judgment. Level 4 governance is about ensuring the former while preventing the latter.

### 4.5.3 Pattern C: Evaluation harnesses

Evaluation harnesses are the least glamorous but most important Level 4 asset pattern. They are the firm's institutional memory about how its AI-enabled assets behave. Without evaluation harnesses, the firm is flying blind. With them, the firm can innovate responsibly.

An evaluation harness is not a one-time test. It is a reusable set of prompts, inputs, and checks designed to probe an asset's behavior under conditions that matter. In consulting, those conditions are rarely clean. Inputs are incomplete. Objectives conflict. Time pressure is high. Ambiguity is normal. A good evaluation harness reflects that reality.

There are three broad categories of tests that should appear in most Level 4 evaluation harnesses. The first category is *baseline behavior*. These tests check that the asset produces the expected structure when given reasonable inputs. They confirm that facts, assumptions, and open questions are separated correctly and that the output matches the intended format. This is necessary but not sufficient.

The second category is *edge cases*. These tests deliberately remove key inputs, introduce ambiguity, or provide contradictory information. The purpose is to see whether the asset responds by surfacing uncertainty or by inventing coherence. For example, an evaluation harness might test a memo generator with no quantitative inputs and check whether the output properly flags missing data rather than fabricating estimates. Or it might test a storyline generator with mutually incompatible objectives and see whether the conflict is made explicit.

The third category is *adversarial or red-team prompts*. These tests try to induce bad behavior. They might explicitly ask the asset to "make reasonable assumptions" without constraints, to "estimate" numbers, or to "add industry benchmarks." The evaluation harness should document whether the asset resists these requests or complies with them. This is not about catching the model out; it is about understanding how easy it is for a user to misuse the asset.

Evaluation results should be written in plain language. They should describe what the asset does well and where it fails. They should not be framed as a scorecard. Consulting firms do not need numerical accuracy metrics to govern Level 4 assets; they need narrative understanding of risk. The evaluation document should also drive design changes. If an asset consistently fails a particular edge case, the response should be to tighten the prompt or add explicit constraints, not to ignore the failure.

Crucially, evaluation harnesses should be saved and versioned alongside the asset. When the asset changes, the harness should be rerun or updated. This creates continuity. It allows the firm to say, "We know how this asset behaves, and we monitor that behavior over time." Without that continuity, evaluation becomes performative.

Finally, evaluation harnesses can be reused across assets. For example, a standard set of red-team prompts that test for invented facts or false authority can be applied to many different prompt libraries and storyline generators. This reuse reduces effort and increases consistency. Over time, the firm develops a shared understanding of its most common AI risks and how to detect them.

### 4.5.4 Pattern D: Training and onboarding kits

The fourth core Level 4 pattern addresses a risk that is often underestimated: human misuse. Even the best-designed asset can be misused if users do not understand its purpose and limits. Training and onboarding kits are therefore not an optional add-on; they are a core asset pattern.

A Level 4 training kit is not a generic "how to use AI" course. It is a targeted set of materials that teach consultants how to use the firm's specific assets safely and effectively. It should cover not only how to invoke an asset, but why the asset is designed the way it is. It should explain the governance logic behind constraints, not just the mechanics.

Effective training kits typically include four components. The first is a conceptual overview that explains the maturity ladder and situates Level 4 assets as scaffolds, not authorities. This reinforces the mental model. The second is a set of worked examples that show good and bad usage. Seeing an example where an asset refuses to invent facts is more powerful than reading a policy statement. The third is a set of exercises that require users to produce and save Asset Bundles or Draft Bundles, reinforcing the habit of documentation. The fourth is a checklist that users can apply before sending anything for review.

Training kits also play a cultural role. They signal what the firm values. If training emphasizes speed and clever prompts, users will optimize for speed. If training emphasizes discipline, explicit assumptions, and review, users will internalize those priorities. Level 4 is where culture and tooling intersect most strongly. Assets teach by example, but training reinforces the lesson.

Onboarding kits should also clarify consequences. Consultants should understand that misuse of assets—such as bypassing review, pasting sensitive client data into prompts, or presenting draft outputs as facts—is not merely a technical error. It is a professional failure. Making this explicit does not require punitive language; it requires clarity. Governance works best when expectations are unambiguous.

From a governance standpoint, training kits themselves should be versioned and owned. As assets evolve, training must evolve with them. A firm that updates assets without updating training creates a gap that will be filled by improvisation. Improvisation is the enemy of Level 4 discipline.

### Why these four patterns are sufficient

It is natural to ask whether these four patterns are enough. What about agent frameworks, automated research, or end-to-end solution builders? In this maturity ladder, the answer is deliberate: those capabilities belong in later levels, if at all. Level 4 is about consolidating what the firm already does well—structured thinking, disciplined drafting, rigorous review—and making it reusable under control.

Prompt libraries, storyline generators, evaluation harnesses, and training kits cover the core of consulting work without pretending to replace judgment. They are also mutually reinforcing. Prompt libraries and storyline generators produce drafts. Evaluation harnesses test them. Training kits teach people how to use them. Together, they create a closed loop of innovation with governance.

Firms that master these patterns will find that they do not need dozens of tools. They need a small number of high-quality assets that are trusted because they are constrained, evaluated, and owned. This is what Level 4 looks like when done well: quiet, disciplined, and effective. It does not announce itself as transformation. It simply makes good consulting easier to do and bad consulting harder to hide.

That is the standard this book sets.

## 4.6 Mini-cases (Level 4): innovation without chaos

Level 4 is where consulting firms face a very specific operational temptation: to confuse "we built something reusable" with "we are ready to scale it everywhere." The firm has tasted leverage. Someone has created a prompt library, a template, or a storyline generator that seems to work. A partner sees a clean draft appear in minutes and naturally asks, "Why aren't we doing this on every engagement?" This is the exact moment when innovation can turn into chaos, because the firm begins to publish reusable assets faster than it can govern them.

These mini-cases are designed to show how Level 4 innovation can be both ambitious and conservative at the same time. Ambitious, because the firm does in fact create reusable internal infrastructure. Conservative, because every asset is bounded, evaluated, owned, and wrapped in controls that keep human judgment in charge. Each case illustrates a class of reusable asset that consulting firms routinely need. Each also illustrates the same core discipline: every asset must encode facts versus assumptions versus open questions; every asset must be transparent about its limits; every asset must produce an audit-ready Asset Bundle; and every asset must be treated as "Not verified" by default.

The mini-cases are deliberately written as firm-wide assets because Level 4 is about reuse across engagements. But they are not written as firm-wide automation, because that is Level 5. In Level 4, the objective is to publish internal assets that reduce variance and improve governance, while still requiring that engagement leadership and reviewers actively supervise outputs. The "innovation without chaos" mantra should be understood as: create leverage, but do not create false authority.

Each case below follows a similar pattern. First, it describes why the asset is attractive and how it creates value. Second, it identifies the specific Level 4 risks that would arise if the asset were published informally. Third, it outlines how the firm packages the asset as an Asset Bundle: scope, exclusions, assumptions, evaluation harness, version control, and ownership. Finally, it shows how the asset is used in practice across multiple engagements without collapsing into a copy-paste machine.

### 4.6.1 Case 1: Firm-wide market entry memo template

A market entry decision memo is one of the most common artifacts in corporate strategy engagements. It appears in some form in growth strategy work, international expansion work, product launch work, and even in M&A contexts where "entering" is equivalent to "acquiring a foothold." The

memo is valuable not because it contains perfect analysis, but because it creates decision hygiene: it forces the team to state the decision, define the options, articulate the tradeoffs, identify risks, and propose next steps. In a consulting firm, this memo is also a training artifact. Junior staff learn how to structure strategy thinking by drafting it.

The firm-wide temptation is obvious: if the memo structure is recurring, why not publish a reusable asset that produces a first draft instantly? This is precisely what Level 4 enables. But the first governance lesson is also obvious: the memo is an argumentative artifact, and LLMs are persuasive narrative engines. If the asset is not constrained, it will begin to invent market facts, competitor dynamics, customer preferences, and plausible but unsupported conclusions. Those inventions will be difficult to spot because they will be embedded inside a structure that already looks like "proper consulting."

The firm therefore publishes a market entry memo template as a Level 4 asset, but only by treating it as a scaffold, not as a market research tool and not as a recommendation engine.

## Asset purpose and scope

The asset's stated purpose is narrow: to create a structured draft memo that can be reviewed, edited, and completed by an engagement team. It is explicitly designed to improve consistency of structure and to surface missing information. It is not designed to provide answers. It is explicitly not designed to generate market sizing, benchmark data, or citations.

Scope is also explicit about lifecycle: it is appropriate for early-stage decision framing and internal alignment drafts, not for final board materials without further work. The asset is labeled as a drafting accelerator for professional consultants, not an authority.

## Explicit exclusions (the heart of safety)

The memo template includes explicit exclusions that appear both in documentation and in the prompt itself. For example:

- The asset must not invent market size, growth rates, competitor shares, customer willingness-to-pay, or regulatory requirements.
- The asset must not cite sources unless the user provides them.
- The asset must not present a final recommendation as fact; recommendations, if requested, must be framed as conditional hypotheses with explicit assumptions.
- The asset must always output open questions when information is missing.

These exclusions are not mere warnings. They are what makes the memo safe to reuse.

## Embedded assumptions and output schema

The asset enforces the consulting governance schema: facts provided, assumptions, open questions, draft output, verification status. This is the primary mechanism for preventing assumption laundering. The memo structure itself is designed to make missing data visible. For example, sections that

typically invite invention—"market size," "competitive landscape," "expected economics"—are written with placeholders. The asset explicitly asks the user to provide verified inputs or else logs open questions such as:

- Provide verified market size and growth (source/date).
- Confirm relevant regulatory constraints for entry.
- Confirm current margin range and cost drivers.

The memo is therefore useful even when data is missing because it clarifies what must be obtained before the memo can become client-ready.

### Evaluation harness for the memo

The firm's evaluation harness for this asset focuses on three predictable failure modes:

1. **False precision.** Does the memo invent numbers or "industry norms" when not provided?
2. **Overconfident recommendation tone.** Does the memo imply certainty when inputs are thin?
3. **Assumption drift.** Does the memo introduce new premises that were not provided by the user?

Test cases include:

- A data-rich market entry scenario where the user provides verified numbers, and the memo must use only those numbers.
- A data-poor scenario where only qualitative context is provided, and the memo must output many open questions and avoid invented detail.
- A negative case where the likely conclusion is "do not enter," to test whether the asset biases toward entry.

The evaluation results are saved in the Asset Bundle with documented failure observations and the changes made to reduce them.

### How reuse works across engagements

Once published, the memo asset is used across multiple engagements, but always through a controlled channel. The asset appears in the firm's registry with a version number. Engagement teams reference the asset version in their Draft Bundle or Workflow Run Bundle. Reviewers can therefore see what scaffold was used.

A key operating rule is enforced: the output of the asset is never client-facing. It is always v0. The engagement team edits to create v1, and a reviewer signs off on v2 before any external distribution. This is not merely a policy; it is embedded in how the asset is taught and how file naming conventions are handled. The memo asset therefore creates leverage without creating chaos: it standardizes structure and governance, while keeping ownership of judgment with humans.

### 4.6.2 Case 2: Cost transformation diagnostic toolkit

Cost transformation work is one of the most common high-pressure consulting engagements. Clients want savings targets, timelines, and initiatives. Teams are often asked to produce early hypotheses quickly, before data is complete. This is exactly the kind of work where AI can help with structure: generating diagnostic question sets, organizing cost levers, drafting workplans, and producing status updates. It is also exactly the kind of work where AI can quietly create liabilities: fabricated savings estimates, generic "best practices," and premature conclusions that become commitments.

The firm therefore builds a Level 4 "cost transformation diagnostic toolkit" as a set of reusable assets. The toolkit is explicitly not a savings estimator. It is a structured diagnostic scaffold designed to help teams ask the right questions, avoid missing categories, and document constraints.

**Toolkit components**

The toolkit is modular, typically including:

- A diagnostic interview guide generator (role-specific question sets).
- A cost category hypothesis map (categories, levers, typical data needs).
- A red-flag checklist (common pitfalls and misleading signals).
- A workplan and PMO cadence template (phases, deliverables, decision gates).
- A weekly update memo template (status, risks, decisions needed, open questions).

Each component is an asset with its own versioning and evaluation, but the toolkit is also packaged as a coherent bundle.

**Scope and exclusions**

The toolkit's scope is diagnostic: hypothesis generation and discovery planning. Its exclusions are explicit:

- It must not generate savings targets or quantify initiative impact without user-provided data and explicit assumptions.
- It must not imply that a lever is feasible without confirming constraints (labor agreements, service levels, regulatory obligations).
- It must not produce "approved initiatives"; it produces hypotheses and required data.

This is crucial because cost transformations are political. A plausible initiative list can quickly become a perceived commitment. The toolkit is designed to resist that drift by using careful language and by forcing open questions.

**Red flags as a governance innovation**

A particularly valuable element of the toolkit is the red-flag checklist, which is itself a Level 4 asset. It encodes lessons learned from prior engagements in a way that is safe to reuse. For example, it flags patterns such as:

- "Savings" that are actually one-time deferrals.
- Initiatives that shift cost rather than reduce it.
- Headcount reductions without workload redesign.
- Vendor renegotiations without understanding contract lock-ins.

The red flags are not presented as truths about the client; they are presented as prompts for scrutiny. This is a good example of governance-first reuse: the firm is not reusing confidential client details; it is reusing disciplined skepticism.

**Evaluation harness**

The evaluation harness tests whether the toolkit:

1. Surfaces constraints rather than assuming them away.
2. Avoids quantification when data is missing.
3. Produces open questions that are actionable (who must answer, what data is needed).
4. Maintains clear distinction between hypothesis and finding.

   Edge cases include:

- A scenario with strict labor constraints.
- A scenario with critical service-level requirements.
- A scenario where the client demands a savings number immediately.

   The goal is to ensure the toolkit does not become a "savings generator," even under pressure.

**Reuse across multiple cost engagements**

In practice, different teams use the toolkit for different industries. The governance wrapper ensures that the toolkit remains generic enough to avoid embedding client-specific details, while still being structured enough to be valuable. Teams are trained to treat the toolkit outputs as starting points for discovery planning. The toolkit's greatest value is not speed; it is completeness. It reduces the chance that a team misses a cost category or fails to ask about a critical constraint. That reduction in blind spots is a defensible form of leverage.

### 4.6.3 Case 3: Capital allocation pre-read generator

Capital allocation work often sits at the most sensitive intersection of strategy, finance, and governance. Pre-read materials for boards or investment committees shape high-stakes decisions. They must be coherent, conservative, and precise about what is known and what is uncertain. This is a domain where a reusable structure can be extremely helpful: it can standardize how options are presented, how tradeoffs are framed, and how risks are documented. It is also a domain where AI misuse can be catastrophic: invented financial metrics, implied recommendations, and overconfident narrative.

For that reason, the firm's Level 4 asset is not a "capital allocation recommender." It is a "pre-read generator" in the narrow sense of producing a controlled narrative structure that forces explicit assumptions and verification planning.

## Why this asset exists

Capital allocation pre-reads often share a common shape:

- Objective and decision required.
- Options under consideration.
- Criteria and tradeoffs.
- Risk assessment and mitigations.
- Implementation considerations.
- Information gaps and verification plan.

Teams spend large amounts of time formatting and rewriting these sections. A reusable asset can reduce that overhead and increase consistency, but only if it is designed to avoid inventing quantitative content.

## Scope and exclusions

The pre-read generator's scope is strictly structural. It can draft sections and suggest questions, but it cannot produce financial numbers unless the user provides them. It cannot generate valuation, discount rates, or projections. It cannot cite sources. It cannot claim compliance with board governance requirements. It must always state verification status as not verified.

The asset explicitly warns that it is unsuitable for final board packs without human review. This is not mere caution; it is the asset's design philosophy. The asset's job is to surface what must be verified and to prevent "answer-shaped" documents from being produced prematurely.

## Assumptions and conditional language

A capital allocation asset must be excellent at conditional reasoning. In many situations, the firm does not yet know the exact return profile or risk distribution. A good asset forces the team to state assumptions explicitly: time horizon, risk tolerance, capital constraints, strategic priorities, and downside scenarios. It also forces the team to separate base case from downside case narrative.

This is where Level 4 discipline becomes visible. The asset is not allowed to collapse uncertainty into a single recommendation. Instead, it produces a structured comparison of options with explicit assumptions. If assumptions are missing, it must log them as open questions. In other words, the asset is designed to be conservative by default.

## Evaluation harness

The evaluation harness tests whether:

1. The asset refuses to invent financial metrics.
2. The asset clearly labels assumptions and separates them from facts.
3. The asset surfaces verification needs and does not imply that work is complete.
4. The tone remains appropriate for governance audiences (clear, cautious, non-speculative).

Test cases include:

- A scenario with complete numbers provided, to ensure the asset uses only the provided numbers.
- A scenario with missing numbers, to ensure the asset produces placeholders and open questions.
- A scenario with politically sensitive tradeoffs (e.g., dividends versus capex), to ensure the asset frames tradeoffs neutrally.

**Reuse without false authority**

The pre-read generator is a classic example of a Level 4 asset that can become falsely authoritative. To prevent this, the firm enforces usage boundaries: the asset can be used for internal drafts and alignment, but final client governance materials require partner review and explicit verification notes. The asset output is never allowed to be labeled "final." This boundary is enforced culturally and operationally through file naming, review gates, and training.

### 4.6.4 Case 4: Operating model redesign playbook

Operating model redesign is one of the most complex types of consulting work because it blends structure and culture. It involves decision rights, processes, roles, governance forums, and incentives. It is also a domain where workshops and interviews are central. Teams need well-designed agendas, interview guides, synthesis templates, and decision logs. These are recurring artifacts, and therefore strong candidates for Level 4 reuse.

At the same time, operating model work is highly contextual. A reusable playbook can easily become a rigid template that pushes a predetermined model. This is a governance and quality risk. The Level 4 playbook must therefore be designed to provide scaffolding while remaining adaptable and explicit about assumptions.

**Playbook components**

A Level 4 operating model redesign playbook typically includes:

- Workshop agenda templates (by phase: current state, future state, transition).
- Interview guides for key stakeholder roles (executives, functional leaders, frontline).
- A decision log template (decisions, owners, dates, rationale, open items).
- A synthesis memo template (themes, tensions, options, risks, next steps).
- Governance reminders embedded in every artifact (facts vs assumptions, open questions).

The playbook is modular, because different engagements require different combinations.

**Scope and exclusions**

The playbook's scope is facilitation and synthesis scaffolding. Its exclusions are explicit:

- It must not prescribe a specific operating model without engagement-specific reasoning and client context.
- It must not claim that a governance forum is "best practice" without stating conditions and assumptions.
- It must not present workshop outputs as decisions unless the decision log records them explicitly.

  The playbook is designed to support professional facilitation rather than to replace it.

**Governance reminders as embedded design**

One of the most important innovations in this playbook is that governance reminders are not placed in a separate "AI policy" document. They are embedded inside the artifacts. For example, the workshop agenda includes a line item: "Record decisions and open questions in the decision log; do not treat discussion as decision." The interview guide includes a section: "Distinguish observations from assumptions; record unknowns." The synthesis memo includes required sections for facts provided, assumptions, and open questions.

This embedding matters because operating model projects generate large amounts of qualitative information. Without disciplined logging, teams can easily drift into narrative certainty. The playbook prevents that by making uncertainty visible and by forcing the team to capture it systematically.

**Evaluation harness**

Evaluation for this playbook focuses less on factual correctness and more on behavioral properties:

1. Does the playbook surface tensions rather than smoothing them over?
2. Does it preserve optionality, or does it push toward a predetermined design?
3. Does it require decision logging and clear ownership?
4. Does it avoid "generic best practice" drift?

  Test cases include different organizational contexts: matrixed global organizations, product-led firms, regulated environments, and high-growth companies. The goal is to ensure the playbook remains adaptable and does not encode a single worldview.

**Reuse across engagements**

The playbook is used widely because it reduces reinvention. Teams can start with a proven workshop sequence and a disciplined interview guide. Reviewers can quickly understand what was done because artifacts are consistent. But the playbook is also designed to require customization: users must specify context, constraints, and objectives; otherwise the playbook outputs more open questions. This prevents "template tyranny," where teams force a client into the shape of the template.

**Cross-case synthesis: what "innovation without chaos" looks like**

Across these four mini-cases, the same Level 4 principles appear repeatedly.

First, assets create value by standardizing *structure*, not by manufacturing content. The market entry memo template standardizes decision hygiene. The cost diagnostic toolkit standardizes discovery planning. The capital allocation pre-read generator standardizes option framing and verification planning. The operating model playbook standardizes facilitation artifacts and decision logging.

Second, assets remain safe because they enforce explicitness: facts provided, assumptions, open questions, and not verified status. This is the core control that prevents AI from becoming a persuasive fiction machine.

Third, assets are governable because they are packaged as Asset Bundles: they have scope boundaries, exclusions, evaluation evidence, version history, and named owners and reviewers. This transforms them from informal templates into managed infrastructure.

Fourth, assets do not remove review. They shift review from rewriting to judgment. Reviewers can focus on the quality of assumptions, the completeness of open questions, and the defensibility of claims, because the asset provides consistent structure.

Finally, and most importantly, these cases show why Level 4 is a governance threshold. Without governance, each asset could easily become a source of institutionalized error: invented market data in the market entry memo, premature savings targets in cost work, implied recommendations in capital allocation materials, and rigid "best practice" operating model prescriptions. With governance, the assets become discipline multipliers: they scale structure, transparency, and auditability.

This is what Level 4 should mean in a consulting firm. Innovation is real, but it is constrained. Leverage is real, but it is supervised. Reuse is real, but it is bounded. The firm becomes better at producing draft work that is consistent, inspectable, and defensible. That is innovation without chaos, and it is the only kind of innovation that can safely serve as a bridge to Level 5, where assets become part of an organizational operating model.

## 4.7   Risks and controls taxonomy (Level 4)

Level 4 is the point in the maturity ladder where the dominant risk is no longer "the model said something wrong." The dominant risk becomes "the firm made something reusable, and that reuse scaled the wrong thing." In Levels 1 through 3, the unit of risk is typically the individual output: a draft email, a memo, a decision record, a multi-step workflow run. Governance can often be enforced through run-level controls: structured schemas, assumption registers, and human review gates. At Level 4, those controls remain necessary, but they are no longer sufficient. The object being governed is now the reusable asset itself, and reuse creates a class of failure modes that do not exist in one-off usage.

This section provides a taxonomy of Level 4 risks and maps them to a minimum control set. The taxonomy is intentionally practical. It is not a theoretical model of AI risk; it is a description

of the specific ways consulting firms tend to fail when they attempt to scale internal templates, playbooks, and prompt libraries. The controls are intentionally conservative. They do not require sophisticated technology. They require discipline, accountability, and lifecycle management.

A useful way to understand this taxonomy is to recognize that Level 4 is a shift from *workflow governance* to *asset governance.* Workflow governance asks: "Was this run reviewed? Were assumptions logged? Was the output labeled not verified?" Asset governance asks: "Should this asset exist? Who owns it? Where is it used? How do we know how it behaves? What changed over time? When do we retire it?" The failure modes unique to reuse arise because assets travel across time and context, and because they carry implied endorsement. These properties make Level 4 powerful and dangerous.

### 4.7.1 Failure modes unique to reuse

1. **Template drift over time.**

Template drift is the most common Level 4 failure mode because it arises naturally from consulting culture. Consulting firms share templates informally. People make small improvements. Partners ask for tone changes. Juniors remove "unnecessary" sections. A prompt that once enforced facts versus assumptions becomes a looser draft generator. None of these changes are malicious. Almost all are well-intentioned. But over time, the asset becomes different from what was originally evaluated and approved.

In an AI-enabled asset world, drift is especially dangerous because small prompt changes can have large behavioral effects. A single sentence like "make reasonable assumptions" can turn a disciplined draft scaffold into an invention engine. A minor change in wording can increase confidence and reduce the visibility of uncertainty. Drift therefore breaks the most important property of governed reuse: predictability.

Drift is also dangerous because it destroys traceability. If a team uses a drifted version of an asset, the firm may believe it is using "the standard tool," but in reality it is using an uncontrolled variant. When something goes wrong, the firm cannot reconstruct what happened because no one can identify which asset version produced the output. In consulting, where credibility depends on being able to explain reasoning and evidence, the inability to reconstruct is itself a governance failure.

2. **Use outside intended context.**

Reuse invites scope creep. A template built for one engagement type will be applied to another. A storyline built for internal alignment will be reused for a governance audience. A diagnostic toolkit built for hypothesis generation will be used as if it were a findings document. This happens because consultants are rewarded for speed and because assets feel like shortcuts.

In traditional template reuse, out-of-context usage is a quality risk. In AI-enabled reuse, it becomes a risk amplifier because the asset can generate persuasive narrative that appears tailored even when the underlying assumptions do not apply. The user may not notice the mismatch. The

reviewer may not notice either, because the output looks structurally "correct." The firm then produces a deliverable that is not wrong in a single obvious sentence but wrong in its overall framing.

Out-of-context use is also a governance risk because it can turn an internal scaffold into a client-facing artifact without the appropriate scrutiny. The firm may not have intended the asset to be used externally, but the boundary is crossed implicitly when a draft is forwarded, pasted into slides, or reused in a workshop. This is why Level 4 boundaries must be explicit and embedded in the asset itself.

3. **Hidden assumptions embedded as defaults.**

Consulting assets inevitably embed assumptions. A template privileges certain questions. A storyline assumes a certain decision logic. A diagnostic toolkit assumes certain cost categories matter. That is not inherently wrong. The failure occurs when these assumptions become hidden defaults and are no longer recognized as assumptions. This is the mechanism of assumption laundering.

In AI-enabled assets, hidden assumptions can be embedded in prompt phrasing ("assume the client has the necessary capabilities"), in default output sections ("industry best practices suggest X"), or in pre-filled narratives. Over time, these defaults become "the way we write," not "the assumptions we made." The firm begins to confuse its own templates with reality.

Hidden assumptions are particularly dangerous in consulting because they shape what teams investigate. If an asset implicitly assumes a certain cause of cost problems, teams will ask questions that confirm it. If an asset implicitly assumes that market entry is attractive, teams will frame options accordingly. The asset becomes a cognitive bias amplifier. Because the output is fluent and structured, the bias is less visible.

A governance-first Level 4 program treats embedded assumptions as design objects that must be explicitly listed, reviewed, and tested. The goal is not to eliminate assumptions; it is to keep them visible.

4. **False sense of validation.**

A reusable asset that looks polished can create the illusion that it has been validated. Teams begin to treat the asset as "approved" and therefore safe. Leaders may assume that because an asset lives in a library, it has been tested. Users may assume that because an asset is "standard," it reflects firm consensus. None of these assumptions are necessarily true.

This false sense of validation is one of the most dangerous Level 4 failure modes because it changes behavior. When teams believe an asset is validated, they become less vigilant. Review becomes lighter. Open questions are ignored. Outputs are treated as closer to final. This is exactly how a firm can begin to outsource discipline without realizing it.

False validation is also an organizational risk because it can become self-reinforcing. The more an asset is used, the more it feels legitimate. But usage is not proof. It is repetition. If the asset is flawed, repetition simply spreads the flaw.

A governance-first firm counters this by insisting that every asset carries an explicit verification status of "Not verified" and by treating evaluation as documentation of failure modes, not as a

certificate of correctness. The asset is never "proven." It is "known," in the sense that its behavior is documented and bounded.

5. **Asset sprawl without ownership.**

As soon as a firm begins to publish reusable assets, it will face asset sprawl. People create multiple versions of similar prompts. Practice groups create parallel toolkits. Teams fork assets for one engagement and never merge improvements back. Over time, the firm accumulates dozens or hundreds of assets, many redundant, many outdated, and many used without clear governance.

Asset sprawl is not merely inefficiency. It is a risk because it destroys consistency and accountability. When multiple assets exist for the same deliverable, teams will choose whichever one produces the most convenient output. This selection bias tends to favor assets that are less constrained and therefore more risky. Sprawl also prevents learning. Failure modes discovered in one asset may not be propagated to others. The firm repeats the same mistakes in parallel.

Ownership is the antidote. An asset without an owner is not an asset; it is a liability. Governance-first Level 4 programs treat ownership as a condition of existence: if no one is willing to own an asset, it should not be published.

### 4.7.2   Minimum control set

The minimum control set below maps directly to the failure modes above. These controls are "non-negotiable" because without them the firm cannot credibly claim that it is innovating responsibly. They are deliberately expressed in plain operational terms. A firm can implement them with simple tools: a controlled repository, a basic registry, a review process, and a set of evaluation test cases. The difficulty is not technical. The difficulty is cultural: enforcing discipline when speed is rewarded.

---

**Checklist**

**Level 4 controls (non-negotiable).**

1. **Asset registry with owners.** The firm maintains a single authoritative list of published assets, each with a named owner and a named reviewer. The registry records purpose, scope, exclusions, current version, last review date, and where the asset is stored. This control addresses asset sprawl and enforces accountability.

2. **Mandatory evaluation before reuse.** No asset is published for reuse without documented evaluation against representative cases and known failure modes. Evaluation does not certify correctness; it documents behavior and limits. The evaluation pack is stored with the asset as part of its Asset Bundle. This control reduces false validation and prevents institutionalizing untested prompts.

3. **Explicit scope and exclusion statements.** Every asset has explicit intended use cases and explicit exclusions. These boundaries appear in documentation and in the asset itself (e.g., within the prompt text and output disclaimers). This control reduces misuse outside intended context and prevents implicit repurposing.

4. **Change control and versioning.** Every asset has a version identifier and a change log. Changes require owner approval and reviewer sign-off for major versions. Informal edits and silent forks are prohibited in the published channel. This control addresses template drift and preserves traceability.

5. **Periodic review and retirement policy.** Assets are reviewed on a set cadence (e.g., semi-annual or annual). Assets that are outdated, redundant, or unsafe are explicitly deprecated and removed from the active library. Retirement is recorded in the registry, and replacement guidance is provided where appropriate. This control limits drift over time and prevents accumulation of orphan assets.

---

## How the controls work together

The controls above are best understood as a closed-loop system rather than as isolated requirements.

The asset registry makes assets visible and ties them to owners. Without visibility, nothing else matters. Mandatory evaluation ensures that assets are published with documented behavior rather than with optimistic anecdote. Scope and exclusion statements ensure that reuse remains bounded and that users cannot plausibly treat assets as general-purpose engines. Change control and versioning ensure that improvements do not become drift and that the firm can reconstruct what happened when outputs are challenged. Periodic review and retirement prevent the library from becoming a graveyard of outdated assets that still circulate informally.

Together, these controls allow the firm to answer the most important governance questions that arise at Level 4:

- What assets exist and who owns them?
- What are they intended to do and not do?

- How do we know how they behave?
- What changed over time and why?
- Are we willing to retire assets that are no longer safe or useful?

A consulting firm that can answer these questions has achieved a real Level 4 capability. It has not merely adopted AI; it has built internal infrastructure with lifecycle governance. That is what makes innovation sustainable.

### A final governance warning

Consulting firms often try to reduce risk through disclaimers. Disclaimers matter, but they are not controls. A disclaimer does not stop drift. It does not prevent misuse. It does not create traceability. It does not assign accountability. Level 4 governance must therefore be operational, not rhetorical. The minimum control set is designed to make governance real: visible assets, owned assets, evaluated assets, bounded assets, versioned assets, and reviewed assets.

If a firm implements these controls consistently, Level 4 becomes a disciplined innovation engine. The firm can publish reusable scaffolds that improve quality and reduce variance. It can learn from failure modes and improve assets over time. It can scale structure without scaling error. If a firm does not implement these controls, Level 4 becomes a factory for institutionalized persuasive error. The difference is not the model. The difference is governance.

## 4.8 Prompt patterns and exercises (copy/paste)

The purpose of Level 4 prompt patterns is not to produce "better outputs" in isolation. The purpose is to help a firm create reusable internal assets that are bounded, inspectable, and governable. In Level 4, a prompt is not a one-off instruction; it is part of an Asset Bundle. That means every prompt pattern must be written as if it will be used by someone who was not present when it was created, under time pressure, with incomplete information, and with a natural temptation to treat fluent text as validated. The prompt patterns in this section are therefore intentionally conservative. They embed scope boundaries. They force assumptions into the open. They require evaluation tests. They label outputs as `"Not verified"` by default.

Below are three concrete examples of how to use the reusable asset creation template. Each example is written as a copy/paste-ready prompt that a firm innovation lab can use to generate an Asset Bundle draft. Each example is designed to be used internally, under supervision, with explicit human review. None of these examples instruct the model to invent facts, benchmarks, or sources. Each includes evaluation tests focused on failure modes typical in consulting reuse: invented facts, assumption drift, false authority, and misuse outside intended scope.

### 4.8.1 Example 1: Market scan memo template (internal draft scaffold)

**Prompt (Copy/Paste)**

```
ROLE: You are assisting in designing a reusable internal consulting asset.
This asset will be used by consulting teams to draft internal memos.
Do not invent facts, market data, benchmarks, citations, or sources.
You must surface assumptions and limits explicitly.
You must design the asset to make missing information visible (open questions),
not to fill gaps.

TASK:
Create a reusable template for a "Market Scan Memo" that helps a consulting
team summarize
an initial view of a market segment for an internal working draft.
The memo is not client-ready and must be labeled as not verified.

REQUIREMENTS:
- The template MUST enforce the separation: facts_provided vs assumptions vs
open_questions.
- The template MUST forbid invented market sizing, growth rates, competitor
shares, or citations.
- The template MUST include an explicit "What we do not yet know" section.
- The template MUST include a "Verification plan" section listing what
must be verified and by whom.
- The template MUST include a short disclaimer that outputs are drafts only.

CONSTRAINTS / EXCLUSIONS (explicitly include these in the asset):
- No external research is performed by the model.
- No citations or data may be added unless provided by the user.
- No numeric estimates may be produced unless provided by the user.

OUTPUT FORMAT:
asset_name:
intended_use:
explicit_exclusions:
assumptions:
template_or_prompt:
evaluation_tests:
verification_status: "Not verified"
```

**How to evaluate this asset (what the innovation lab checks)**

This asset is successful if it reliably produces a disciplined *structure* for a market scan memo without smuggling in invented market facts. The innovation lab should pay special attention to sections

that usually attract hallucination (market size, growth, competitive landscape). A good template refuses to fabricate and instead produces open questions.

Suggested evaluation lens:

- **No new facts test:** Provide only qualitative notes and confirm the output does not introduce numbers.
- **Ambiguity test:** Provide conflicting notes (e.g., "market is growing" and "demand is declining") and confirm the output surfaces the conflict instead of resolving it.
- **Scope misuse test:** Ask the asset to "estimate market size" and confirm it refuses and logs an open question instead.
- **Tone test:** Confirm the output does not read as a client recommendation; it remains an internal draft with explicit verification status.

## 4.8.2 Example 2: PMO weekly update template (status clarity + decision hygiene)

**Prompt (Copy/Paste)**

```
ROLE: You are assisting in designing a reusable internal consulting asset.
This asset will be used to generate a weekly PMO update draft for ongoing engagement.
Do not invent facts, dates, owners, milestones, risks, or metrics.
You must surface assumptions and limits explicitly.
Your job is to produce a disciplined status structure that makes gaps visible.


TASK:
Create a reusable template for a "PMO Weekly Update" used by consulting teams.
The update should support decision hygiene:
what changed, what is blocked, and what decisions are needed.



The output must be a draft and must be reviewed before client sharing.


REQUIREMENTS:
- "Separation": facts_provided vs assumptions vs open_questions.
- "Decisions needed this week" section.
- "Risks and mitigations" section, but MUST NOT invent risks; if not provided,
  it must request them as open questions.
- "Actions / owners / due dates" table format, but MUST NOT create owners
  or due dates unless provided by the user.
- "Client communications readiness" line stating verification_status.

CONSTRAINTS / EXCLUSIONS (explicitly include these in the asset):
- No new milestones, dates, or owners may be invented.
- If inputs are missing, output MUST add open_questions rather than guessing.
- No claims about progress percentages unless provided by the user.

OUTPUT FORMAT:
asset_name:
intended_use:
explicit_exclusions:
assumptions:
template_or_prompt:
evaluation_tests:
verification_status: "Not verified"
```

**How to evaluate this asset (what the innovation lab checks)**

This asset is successful if it reduces variance in PMO updates without accidentally creating fictional project reality. Weekly updates are particularly prone to "helpful completion," where a model fills in owners, dates, or status language to make the update look tidy. The innovation lab should test for this explicitly.

Suggested evaluation lens:

- **Missing owners test:** Provide actions without owners and confirm the output asks for owners rather than inventing them.
- **Missing dates test:** Provide milestones without dates and confirm the output logs open questions and preserves placeholders.
- **Overconfidence test:** Provide ambiguous progress notes and confirm the output avoids false precision and clarifies uncertainty.
- **Structure stability test:** Run the asset on three different projects and confirm the structure remains stable and review-friendly.

### 4.8.3 Example 3: Workshop agenda + decision record kit (operating model redesign)

---

**Prompt (Copy/Paste)**

```
ROLE: You are assisting in designing a reusable internal consulting asset.
This asset will be used to plan workshops and to capture outcomes as a decision record.
Do not invent facts about the client, stakeholders, governance forums, or decisions.
You must surface assumptions and limits explicitly.
You must design the asset to distinguish discussion from decisions.


TASK:
Create a reusable template kit for
"Operating Model Workshop Agenda + Decision Record".


The kit should include:


(1) a workshop agenda template with timeboxes and objectives,
(2) a decision record template to be completed during or
immediately after the workshop.


REQUIREMENTS:
- The templates MUST enforce the separation: facts_provided vs assumptions
vs open_questions.
- The agenda MUST include explicit points where decisions are intended,
and instruct users to capture them
   in the decision record (not in narrative notes).
- The decision record MUST separate: decisions
made vs actions agreed vs open issues.
- The kit MUST include governance reminders:
facts are not assumptions; outputs are drafts; verification required.
- The kit MUST avoid "best practice" prescriptions unless
explicitly provided by the user as a chosen principle.


CONSTRAINTS / EXCLUSIONS (explicitly include these in the asset):
- Do not invent attendees, roles, decision rights, or governance structures.
- Do not fabricate consensus; if alignment is not provided, log it as an open question.
- Do not create a future-state operating model
recommendation; this is facilitation scaffolding only.


OUTPUT FORMAT:
asset_name:
intended_use:
explicit_exclusions:
assumptions:
template_or_prompt:
evaluation_tests:
verification_status: "Not verified"
```

**How to evaluate this asset (what the innovation lab checks)**

This asset is successful if it prevents one of the most common consulting failures: treating workshop discussion as a decision. It should also prevent narrative smoothing, where tension is rewritten into alignment. The innovation lab should evaluate not only whether the templates look professional, but whether they enforce disciplined capture of uncertainty, disagreement, and open issues.

Suggested evaluation lens:

- **No-attendees test:** Provide no participant list and confirm the output requests attendees rather than inventing them.
- **No-decisions test:** Provide notes that contain discussion but no decisions and confirm the decision record clearly indicates "no decision made" rather than implying one.
- **Conflict test:** Provide notes indicating disagreement and confirm the output preserves the disagreement and records it as an open issue.
- **Misuse test:** Ask the kit to "recommend the target operating model" and confirm it refuses and redirects to open questions and facilitation outputs.

### 4.8.4 Exercise set (firm innovation lab)

The following exercises are designed to train a firm innovation lab to build, test, and review Level 4 assets without creating uncontrolled sprawl. The exercises assume that every asset produced will be packaged as an Asset Bundle and stored in a controlled repository with a named owner and reviewer. The goal is not creativity; the goal is disciplined reuse.

**Exercise 1: Build one asset, then prove its boundaries**

Choose one of the three example assets above and produce a first Asset Bundle draft. Then conduct a boundary proof exercise:

1. Run the asset with minimal inputs and record whether it invents content.
2. Ask the asset explicitly to do something out of scope (e.g., add benchmarks, estimate savings, cite sources).
3. Verify that the asset refuses and logs an open question instead.
4. If it fails, revise the prompt and rerun until refusal behavior is reliable.

Deliverable: a short evaluation note documenting the misuse attempts and the resulting changes to the asset.

**Exercise 2: Drift simulation (what happens when helpful edits accumulate)**

Take a published asset draft and intentionally introduce three small "helpful" edits (for example: "make reasonable assumptions," "add typical benchmarks," "fill in owners and dates"). Then rerun the evaluation harness and document how behavior changes.

1. Identify which edit caused the most serious governance regression.

2. Write a change-control rule that would prevent that edit from being merged silently.

3. Add that rule to the asset's explicit exclusions.

Deliverable: a drift report that demonstrates why versioning and review gates are non-negotiable.

### Exercise 3: Reviewer simulation (partner-quality scrutiny in 20 minutes)

Assign one person to act as asset owner and one to act as reviewer. The reviewer has 20 minutes to inspect the Asset Bundle and answer:

- Is the scope sufficiently narrow?
- Are exclusions explicit and embedded in the prompt?
- Are assumptions visible or hidden?
- Does the evaluation harness test realistic failure modes?
- Could a junior misuse this asset easily?

If the reviewer cannot answer these questions from the bundle, the bundle is incomplete. Deliverable: reviewer notes and required edits before internal publication.

### Exercise 4: Library hygiene (preventing sprawl)

As a group, propose five new assets the firm "should build." Then apply a hygiene filter:

1. Identify which proposals are duplicates of existing assets.
2. Identify which proposals are too broad (attempting to "solve consulting").
3. Choose only two assets to proceed with and write explicit reasons why the others are deferred.

Deliverable: an innovation backlog with clear selection rationale, demonstrating that governance includes saying "no."

### Exercise 5: Publish-and-trace drill (auditability of reuse)

Simulate publication of an asset (assign version v1.0) and then simulate three engagement teams using it. Each team produces a Draft Bundle that references the asset name and version.

1. Confirm that each Draft Bundle can be traced back to the correct asset version.
2. Confirm that changes made by teams are captured as human edits, not as silent asset modifications.
3. Confirm that any discovered failure modes are fed back to the owner for potential version updates.

Deliverable: three Draft Bundles + a short traceability note showing how the firm reconstructs usage.

**Exercise 6: Retirement decision (the discipline to deprecate)**

Pick an asset and imagine it has become obsolete or risky (e.g., repeated misuse, new failure mode, changed firm policy). Write a retirement plan:

1. Define retirement criteria (what triggers deprecation).
2. Draft a deprecation notice for the registry.
3. Define replacement guidance (if any).
4. Define how to prevent continued informal circulation.

Deliverable: a retirement policy snippet that can be reused across the asset library.

### What "good" looks like after these exercises

After completing these exercises, an innovation lab should be able to demonstrate Level 4 maturity in practical terms. They should be able to produce Asset Bundles that are small, clear, and governable. They should be able to show evaluation evidence focused on failure modes rather than success stories. They should be able to enforce change control to prevent drift. They should be able to trace asset usage across engagements. And they should be comfortable retiring assets when governance requires it.

These are not technical achievements. They are professional discipline achievements. They represent the consulting firm learning to treat AI-enabled assets as infrastructure: powerful, reusable, and therefore deserving of lifecycle governance.

## 4.9 Conclusion and transition to Level 5 (Organizations)

### 4.9.1 Summary of main takeaways

Level 4 represents a decisive shift in how a consulting firm relates to AI. In earlier levels, AI is something teams *use*. In Level 4, AI becomes something the firm *builds around*. The central achievement of this level is not speed, sophistication, or novelty. It is the disciplined creation of reusable internal assets that improve consistency and governance without eroding professional judgment.

The defining insight of Level 4 is simple but easy to ignore: reuse changes the risk profile. A one-off prompt that fails affects a single deliverable. A reusable asset that fails affects the firm. It shapes default behavior, influences how problems are framed, and carries implied endorsement. This is why Level 4 is not primarily an innovation challenge; it is a governance challenge disguised as innovation.

Throughout this chapter, the same pattern has appeared repeatedly. Reusable assets are valuable when they standardize *structure*, not when they manufacture content. Prompt libraries help teams draft consistently without inventing facts. Storyline generators clarify logic without pretending to supply evidence. Diagnostic toolkits surface hypotheses and constraints rather than premature

conclusions. Evaluation harnesses document failure modes instead of certifying correctness. Training kits teach discipline rather than shortcuts. In every case, the asset is a scaffold, not an authority.

The biggest risk at Level 4 is therefore not model error in the narrow sense. It is institutionalizing error. A hidden assumption becomes a default. A persuasive tone becomes a substitute for verification. A "standard template" becomes an excuse to stop thinking. These risks do not arise because consultants are careless. They arise because consulting firms are under pressure to move quickly and because reusable assets feel trustworthy once they are published. Without governance, the firm confuses repetition with validation.

The strongest control at Level 4 is ownership reinforced by evaluation. Ownership ensures that every asset has someone accountable for its purpose, boundaries, and evolution. Evaluation ensures that the firm understands how the asset behaves under realistic conditions, including failure modes. Together, ownership and evaluation turn innovation into something the firm can stand behind. They make it possible to say, credibly, "We know what this asset does, we know where it fails, and we know who is responsible for it."

Just as importantly, Level 4 reframes what "progress" looks like. Progress is not a growing library of clever prompts. Progress is a shrinking gap between how the firm intends to work and how it actually works under pressure. When assets are governed well, they reduce variance without reducing accountability. They make it easier for junior teams to produce reviewable drafts. They make review more focused on judgment rather than formatting. They make it easier for leadership to understand what methods are actually being used across engagements.

In that sense, Level 4 is a maturity test. Firms that treat AI assets casually will experience sprawl, drift, and false authority. Firms that treat AI assets as internal infrastructure will discover that a small number of well-governed assets can outperform a large number of improvised tools. The difference is not technical capability. It is discipline.

### 4.9.2 What comes next (preview of Level 5)

Level 4 deliberately stops short of firm-wide automation. Assets exist. They are reusable. They are governed. But they are still optional tools used by teams within engagements. Human judgment is local. Review is engagement-specific. Governance, while real, is still largely asset-centric.

Level 5 changes the unit of analysis again.

In Level 5, AI-enabled assets are no longer just tools in a library; they become components of the firm's operating model. Governance moves from asset-level discipline to organization-level systems. The firm begins to answer questions that cannot be answered at Level 4 alone: Who is allowed to use which assets, in which contexts, and with what approvals? How are AI-influenced outputs routed, reviewed, and archived at scale? How does the firm monitor usage patterns and detect emerging risks? How does it demonstrate auditability, not just internally but to regulators, clients, and boards?

Where Level 4 is about publishing assets responsibly, Level 5 is about embedding those assets into controlled workflows. Intake processes define what data can be used and how it must be redacted.

Approval routing ensures that sensitive outputs receive the right level of review. Monitoring systems track which assets are used where and how often. Recordkeeping systems allow the firm to reconstruct decisions long after the engagement ends. Governance becomes systemic rather than artisanal.

This is also where organizational accountability becomes explicit. In Level 4, accountability is anchored in named asset owners and reviewers. In Level 5, accountability expands to committees, policies, escalation paths, and firm-wide standards. Questions of independence, confidentiality, conflicts, and risk appetite are no longer handled implicitly by individual partners; they are encoded into how the firm operates. AI is no longer an innovation initiative. It is part of how the firm delivers work.

Crucially, Level 5 does not replace Level 4 discipline; it depends on it. A firm cannot build an organizational operating model on top of unmanaged assets. Asset registries become inputs to firm-wide asset governance. Evaluation harnesses become part of quality assurance. Version control becomes part of compliance. Traceability becomes a management requirement rather than a best practice. Level 5 is only possible if Level 4 has been done properly.

The transition from Level 4 to Level 5 should therefore feel less like a technological leap and more like an organizational commitment. The firm decides that AI-enabled assets are important enough to deserve the same rigor it applies to other core elements of delivery: methodology, quality control, and risk management. This is where AI governance stops being a chapter topic and becomes part of how the firm runs itself.

---

**Artifact (Save This)**

**Level 4 exit criteria (ready to move to Level 5).**

1. The firm maintains an AI asset registry with named owners and reviewers, and leadership treats it as authoritative.
2. Reusable assets have documented evaluation results focused on failure modes, not just success cases.
3. Change control and versioning are enforced, and silent drift is no longer tolerated in published assets.
4. Leadership can trace asset usage across engagements and can reconstruct which assets influenced which deliverables.

---

Reaching these exit criteria does not mean the firm has "solved" AI. It means the firm has earned the right to move from innovation to organization. It has demonstrated that it can reuse without losing control, standardize without diluting judgment, and innovate without institutionalizing error.

Level 5 will build on this foundation. It will show how consulting firms turn governed assets into governed systems, and how governance itself becomes a competitive advantage. But that step is only credible if Level 4 has been taken seriously. In this maturity ladder, there are no shortcuts. Every increase in capability demands a corresponding increase in responsibility. Level 4 is where

the firm proves it understands that rule.

# Bibliography

[1] National Institute of Standards and Technology (NIST). *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*. NIST AI 100-1, U.S. Department of Commerce, January 2023.

[2] International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC). *ISO/IEC 42001:2023 — Information technology — Artificial intelligence — Management system*. ISO/IEC Standard, December 2023.

[3] Board of Governors of the Federal Reserve System. *Supervisory Guidance on Model Risk Management (SR 11-7)*. Supervisory Letter SR 11-7 and attachment, April 4, 2011.

[4] Office of the Comptroller of the Currency (OCC). *Sound Practices for Model Risk Management: Supervisory Guidance on Model Risk Management*. OCC Bulletin 2011-12, April 4, 2011.

[5] National Institute of Standards and Technology (NIST). *Secure Software Development Framework (SSDF) Version 1.1: Recommendations for Mitigating the Risk of Software Vulnerabilities*. NIST Special Publication 800-218, February 2022.

[6] OECD. *Recommendation of the Council on Artificial Intelligence*. OECD Legal Instrument OECD/LEGAL/0449, May 22, 2019.

[7] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model Cards for Model Reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency (FAccT)*, 2019.

[8] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. Datasheets for Datasets. *Communications of the ACM*, 64(12), 2021.

# Chapter 5

# Organizations

**Abstract.** This chapter defines Level 5 (Organizations) maturity for AI consulting companies: the point at which generative AI is no longer a set of helpful tools, but a governed delivery operating model. The chapter's stance is conservative. AI accelerates drafting and structured analysis, but it does not replace professional judgment, and all outputs remain *Verification status: Not verified. Human review required.* until verified by humans. Level 5 is characterized by firm-wide controls that make work reconstructable: an AI asset registry with owners and permitted-use statements; standardized intake and sensitivity classification; stage-gated workflows with separation of duties; mandatory structured outputs that separate facts, assumptions, and open questions; and evidence bundles that link every material claim to a source or an explicit assumption with a test plan. The chapter also introduces the Consulting Workpaper Bundle (CWB) as the minimum recordkeeping unit for AI-assisted engagements, and it specifies evaluation harnesses, monitoring, incident response, and change control to prevent template drift and systemic error. Four recurring consulting mini-cases illustrate how governance scales safe reuse across teams and clients.

---

**Artifact (Save This)**

**Scope and stance (read before using).** This chapter is written for consulting-company leaders (partners, principals, practice leaders, transformation directors, COO/CFO office, knowledge management, risk/compliance, and delivery leaders) with minimal AI background. It focuses on **Level 5 maturity: organization-level AI operating models**, where workflows, controls, recordkeeping, and supervision are embedded into the firm's delivery system. The operating posture is conservative: **AI does not replace professional judgment**, **facts must be verified by humans**, and **every client-facing output remains *Verification status: Not verified. Human review required.* until validated**. This is educational and workflow-focused content; it is not legal, financial, or regulatory advice.

## 5.1 Chapter overview: where Level 5 fits in the maturity ladder

**Why this chapter exists.** By Level 5, the consulting firm has already learned the easy lesson: a well-prompted model can draft faster, restructure slides, and produce plausible prose on demand. That lesson matters, but it is not the organizational problem that determines whether AI becomes an asset or a liability. The Level 5 problem is scale. When hundreds (or thousands) of drafts, analyses, and client-facing deliverables are produced under time pressure across multiple teams, the core question becomes: can the firm reliably produce defensible work at scale, with consistent quality, clear accountability, and reconstructable evidence?

Consulting firms sell judgment, credibility, and trust. A client may disagree with a recommendation, but they should not be able to credibly say, "your work was careless, unverifiable, or procedurally unserious." Ungoverned AI puts that trust at risk in precisely the ways that matter most in consulting: it can create fluent but ungrounded claims; it can conceal uncertainty behind confident wording; it can accelerate the spread of template drift; it can import external content that is not permitted; and it can blur the line between what the client provided, what the team assumed, and what the team verified. The problem is not merely that models can be wrong. The deeper problem is that models can be wrong *persuasively*, and when persuasive error is multiplied across teams through reuse, the organization does not merely produce one defective output. It industrializes the defect.

Level 5 exists to prevent the firm from becoming a factory for confident nonsense. It does so by making governance part of delivery. That phrase is not a slogan; it is a design requirement. The firm does not rely on individual discipline or partner intuition alone. Instead, it builds a delivery operating model in which the correct process is the default, and deviations are visible. At Level 5, the unit of value is not a prompt, and it is not even a workflow. The unit of value is a governed *system* that turns client inputs into deliverables through defined stages, defined roles, defined approvals, and defined artifacts, producing evidence as a byproduct of doing the work.

This is why Level 5 is the natural end of a governance-first maturity ladder. Levels 1–4 can be adopted by motivated individuals and small teams. Level 1 teaches disciplined drafting. Level 2 makes reasoning inspectable by forcing structure: issue trees, alternatives, tradeoffs, and explicit assumptions. Level 3 introduces multi-step workflows with human checkpoints, separation of duties, and immutable logs. Level 4 turns those workflows into reusable assets with evaluation harnesses, controlled releases, and supervision frameworks. But none of those levels guarantees firm-wide reliability on their own. They increase capability and introduce local controls. They do not yet answer the organization-level question: how does a consulting company ensure that AI-assisted delivery is consistent, auditable, and safe across offices, service lines, partners, and client contexts?

Level 5 answers by treating AI as an enterprise delivery component. The firm manages it like any other production system: with ownership, permitted use, change control, monitoring, and incident response. The firm defines what "good" looks like operationally, not rhetorically. "Good" does not mean the deck looks polished. "Good" means that every material claim can be traced to evidence or to an explicit assumption with a test plan; that confidential information is handled under

minimum-necessary rules; that outputs are labeled *Verification status: Not verified. Human review required.* until verified; that review and approvals are recorded; and that the firm can reconstruct what was produced, when, by whom, using which assets and versions, and why it was reasonable to present it.

Notice the shift in posture. Level 5 does not promise to reduce risk to zero. It promises to make risk manageable, visible, and bounded. It recognizes that consulting work is inherently uncertain: strategy involves forward-looking judgments and incomplete information. The goal is not to eliminate assumptions; it is to keep assumptions from masquerading as facts, and to ensure that the firm knows where the uncertainty is concentrated. In practice, Level 5 changes how the firm thinks about AI-supported work: not as outputs that must be "trusted," but as drafts and intermediate artifacts that must be governed.

This chapter therefore exists for leaders, not for hobbyists. It is written for partners who sign engagement letters, for delivery leaders who run multi-team programs, for knowledge management leaders who publish firm templates, for risk and compliance leaders who worry about confidentiality and reputational damage, and for operations and technology leaders who must implement controls that are real rather than aspirational. It treats the consulting company as the primary unit of analysis. Individuals still matter, but Level 5 assumes that individual discipline is insufficient. The firm must build a system that protects clients, protects the firm, and protects the integrity of professional judgment.

**Objectives.**

1. **Define Level 5 as an organization-level operating model (not a tool, not a prompt library).** At Level 5, maturity is not measured by which model the firm uses, or how many templates it has, or how frequently employees prompt an assistant. Maturity is measured by whether the firm can specify and enforce a delivery lifecycle that produces reconstructable evidence. This objective establishes the core mental model: Level 5 is an operating system for AI-assisted consulting delivery, with defined stages, controls, and artifacts.

2. **Provide a concrete blueprint for stage gates, roles, approvals, and recordkeeping.** A governance-first firm cannot depend on informal "someone probably checked" norms. This objective specifies the minimum architecture of a stage-gated workflow, including separation of duties (producer, reviewer, approver), required attestations, and the minimum recordkeeping unit that must exist for any reliance-bearing client artifact. The chapter emphasizes that controls must be embedded in process, not appended after the fact.

3. **Teach how to manage AI assets like any other production system: versioning, evaluation, monitoring, and deprecation.** At Level 5, a prompt template is not a personal convenience; it is a firm asset that can spread errors. This objective introduces asset governance: an AI asset registry, ownership, permitted-use statements, input/output schemas, evaluation harnesses and regression tests, release notes, and deprecation policies. It treats templates, playbooks, and model configurations as controlled artifacts with lifecycle management.

4. **Show how to scale reuse without scaling error: templates, playbooks, and test**

**harnesses with change control.** Reuse is where consulting firms gain leverage, and it is also where AI creates the highest systemic risk. This objective focuses on the discipline of reuse: how to standardize deliverables without standardizing hallucinations; how to ensure that what is reused is structure and governance rather than unverified content; and how to require that every change to reusable assets is tested and approved before it is pushed firm-wide.

5. **Produce a minimum viable Consulting Workpaper Bundle (CWB) standard for AI-assisted delivery.** A mature firm does not debate recordkeeping ad hoc on every engagement. This objective defines the CWB as the minimum set of artifacts that makes an AI-assisted deliverable defensible: facts provided, assumptions register, open questions, verification plan, evidence map, model/workflow lineage, outputs and edits, reviewer notes, approvals, and risk flags. The CWB is not paperwork for its own sake; it is the mechanism that preserves accountability and enables reconstruction.

**The five-level maturity ladder (preview).** The maturity ladder in this book is deliberately simple, because simplicity is enforceable. Level 1 (**Chatbots**) treats AI as a drafting assistant: accelerate writing while enforcing a conservative verification posture and clear separation of facts, assumptions, and open questions. Level 2 (**Reasoners**) adds explicit reasoning structures: issue trees, alternatives, tradeoffs, and assumption registers that make thinking inspectable and reduce decision laundering. Level 3 (**Agents**) introduces multi-step workflows with human checkpoints and immutable logs: the process becomes the artifact, and separation of duties is enforced through stage gates. Level 4 (**Innovators**) turns workflows into reusable assets: standardized templates, playbooks, evaluation harnesses, and controlled releases so that reuse is safe rather than contagious. Level 5 (**Organizations**) completes the shift from individual tool use to firm-level operating model: policy, workflow, evaluation, monitoring, incident response, and recordkeeping are integrated into delivery so that the firm can scale AI assistance without scaling unverified claims, confidentiality failures, or reputational damage.

The ladder's governing principle is non-negotiable: as capability increases, risk increases, and controls must increase in parallel. The purpose is not to make the model "smarter." The purpose is to make the firm *more defensible*. In a consulting context, defensibility means that the firm can show its work. It can demonstrate how inputs were handled, how claims were formed, what was verified, what remained uncertain, who approved what, and what evidence supports the client-facing narrative. That is the Level 5 standard, and it is why this chapter sits at the top of the ladder. It is where the firm stops asking whether AI is useful and starts proving that AI-assisted delivery is governed.

## 5.2 Mental model: what Level 5 "Organizations" actually means

**Useful abstraction.** The most useful way to understand Level 5 is to treat it as a *delivery operating system*. Not a tool, not a chat interface, not a library of prompts, and not a procurement decision. An operating system is the layer that determines what is allowed to run, how work is scheduled,

what resources can be touched, what gets logged, what fails safely, and what evidence remains after the work is done. Level 5 applies that same logic to consulting delivery: the firm defines standard inputs, standard outputs, and standard checkpoints, and it treats AI as a controlled component embedded in a supervised pipeline.

This is a deliberately unromantic framing. It avoids the two extremes that distort implementation. The first extreme is the "AI will replace consulting" fantasy, which confuses fluent language with professional judgment and underestimates governance. The second extreme is the "AI is just like Excel" analogy, which treats AI as a neutral productivity tool and ignores that it can generate plausible but incorrect content, and can do so in ways that are hard to detect under time pressure. Level 5 rejects both. It treats AI as a powerful drafting and structuring engine that is valuable precisely because it is fast and flexible, but risky precisely because it can be persuasive without being grounded.

In practice, the delivery operating system abstraction implies five commitments.

First, **the firm standardizes inputs**. At Level 5, teams do not begin with a blank page and an unstructured prompt. They begin with an intake object: a controlled representation of what is known, what is confidential, what is in scope, what sources are permitted, what decisions are being supported, and what verification is required before any reliance-bearing statement may leave the firm. Standardized input does not mean all engagements are identical. It means that every engagement begins with the same governance questions: What are the facts provided? Who owns them? What is missing? What assumptions are we introducing? What evidence will we need? What is the sensitivity class of the data? Which tools are permitted for this class? Which outputs are client-facing? What approvals are required?

Second, **the firm standardizes outputs**. Level 5 does not optimize for "beautiful prose." It optimizes for traceability and reviewability. Outputs are structured by default: facts versus assumptions versus open questions; draft narrative; verification plan; evidence map; risk flags; and handoff notes. This structure is not aesthetic. It is a safety rail. It forces the organization to keep uncertainty visible and to prevent "facts" from silently changing as drafts evolve. It also makes downstream review feasible. A reviewer cannot efficiently validate a ten-page narrative that blends sourced claims with guesses. They can validate a structured packet where the claims are enumerated, labeled, and mapped to evidence.

Third, **the firm standardizes checkpoints**. In consulting, the real product is not the slide; it is the judgment behind it. Level 5 therefore puts stage gates where judgment can be corrupted: scoping and sensitivity classification before data is ingested; claims verification before quantitative statements are presented; peer review before conclusions are stabilized; partner approval before external delivery; and recordkeeping retention immediately after delivery. Checkpoints are not optional reminders. They are enforced transitions that block progress when evidence, approvals, or governance artifacts are missing. The firm is not "more mature" because it tells people to be careful. It is more mature because it builds workflows that fail closed rather than fail open.

Fourth, **the firm standardizes accountability**. AI makes it easy for responsibility to dissolve

into the system: "the model said," "the template produced," "the platform generated." Level 5 exists to prevent that diffusion. Every artifact has owners: an engagement owner, a workstream owner, a reviewer, and an approver. Every assumption has an owner. Every reusable template has an owner. Every model configuration has an owner. Ownership is not ceremonial. It is what makes it possible to ask: who is accountable for verifying this claim, for approving this deliverable, for maintaining this template, for monitoring this asset in production?

Fifth, **the firm standardizes evidence**. At Level 5, evidence is not something you remember to attach after a client challenge. Evidence is a byproduct of doing the work. The operating system produces logs, hashes, version identifiers, reviewer notes, and retention bundles automatically. The goal is reconstructability: the ability to reconstruct what was produced, with what inputs, under what constraints, by which workflow version, using which assets, and with which approvals. This is not merely risk theater. In a consulting context, reconstructability protects the client and protects the firm. It makes it possible to correct errors, learn from incidents, and demonstrate professionalism under scrutiny.

The delivery operating system abstraction also clarifies what Level 5 is *not.* It is not a claim that AI outputs are more accurate than humans. It is not a claim that governance can be outsourced to a vendor. It is not a claim that compliance is solved by turning on enterprise SSO. It is not a claim that the firm can do more speculative work without disclosure. It is, instead, a claim that the firm can operationalize discipline so that speed does not erase defensibility.

To make this concrete, imagine two firms producing a market entry deck. Both firms may use the same model. Both may have talented partners. The Level 5 firm differs in one decisive way: it can demonstrate how its claims were formed, what was verified, what was assumed, what remains uncertain, and how the client should interpret the outputs. That is the difference between a firm that uses AI and a firm that has governed AI.

**Dangerous misconception.** The first dangerous misconception is the most common: "If we have a firm-wide AI platform, we are governed." This misconception is attractive because it makes governance feel like procurement. Buy an enterprise tool, control access, and declare victory. But governance is not an access control problem; it is an accountability and evidence problem. A platform can restrict who can log in. It cannot, by itself, ensure that a claim is sourced, that an assumption is labeled, that a review occurred, that a stage gate was respected, or that a deliverable can be reconstructed. Those outcomes require workflow design and operational discipline.

Access controls are necessary, but they are not sufficient. They are the perimeter fence, not the internal discipline. A firm can have perfect SSO and still ship decks with invented benchmarks. It can have a secure environment and still suffer decision laundering: using AI-generated narrative to justify a pre-decided conclusion while pretending the process was analytical. It can lock down data and still fail to retain evidence, making it impossible to respond credibly when a client asks, "Where did this number come from?"

The second dangerous misconception is subtler: "Evaluation is a one-time model selection." Firms often treat evaluation like vendor due diligence. They test a model, choose a provider, and

assume the risk is controlled. Level 5 insists on a different reality: evaluation is continuous because the system changes continuously. Models change. Prompt templates evolve. Teams adapt. New use cases appear. The distribution of inputs shifts. What was safe last quarter can become unsafe next quarter because a template drifted, because a model update changed behavior, because new consultants copied a pattern without understanding constraints, or because a new client context introduced sensitive data that was not present in the original design.

In other words, evaluation is not a gate at the beginning of the journey. It is a guardrail throughout the journey. The firm must evaluate before releasing assets, after modifying assets, and periodically during usage. It must monitor outputs for drift: changes in tone that creep toward overconfidence, changes in phrasing that weaken disclaimers, changes in the frequency of unverified claims, or changes in the prevalence of sensitive data leakage. Without monitoring, the firm will eventually confuse "we used to be careful" with "we are careful."

The third misconception is strategic: "Level 5 means we automate judgment." The temptation at Level 5 is to treat the organization operating model as a justification for automation: if we have controls and logs, why not let the system decide? That is precisely backwards. The more mature the system, the more clearly it should distinguish automation of *process* from ownership of *judgment*. Consulting value is grounded in judgment under uncertainty. The firm can automate formatting, summarization, and the production of structured drafts. It cannot automate responsibility. If Level 5 becomes a pathway to autonomy, it will eventually become a pathway to institutionalized error. The firm will be fast, consistent, and wrong, which is the worst possible combination.

The fourth misconception concerns "professionalism by formatting." AI can produce outputs that look like consulting deliverables: issue trees, SWOTs, business cases, and polished slide narratives. This can create false rigor. When something is formatted like analysis, it can be treated as analysis even when it is not grounded. Level 5 must treat "consulting-shaped text" as a risk, not as a success criterion. The firm should assume that persuasive formatting increases the probability of uncritical acceptance. Therefore, governance must intensify as outputs become more polished: stronger verification, clearer assumptions, and more explicit limitations.

The fifth misconception is cultural: "Governance is a compliance layer that slows delivery." In reality, governance is a delivery enabler when it is designed correctly. A firm without Level 5 governance eventually slows down anyway, but for worse reasons: rework from errors, partner escalations, client disputes, and internal distrust of outputs. Level 5 is meant to create speed with discipline. It allows reuse because assets are tested. It allows delegation because roles are clear. It allows scale because evidence is produced automatically rather than reconstructed under stress.

A useful test is to ask: can the firm explain its AI-assisted process to a skeptical client and still sound credible? If the explanation is "we use an enterprise model" or "our platform is secure," the firm is not governed. If the explanation is "we have stage gates, we separate facts from assumptions, we verify material claims against approved sources, we retain workpaper bundles, and a partner signs off before delivery," the firm is approaching Level 5.

**Definition of "good" Level 5 output.** At Level 5, "good output" is not impressive text.

It is not creative language. It is not the illusion of certainty. Good output is an artifact that a professional can rely on *after verification*, and that the firm can defend under scrutiny. That means five things.

First, good output is **correct or explicitly unverified**. Consulting deliverables often involve forward-looking scenarios and qualitative judgments; correctness does not mean certainty. It means that factual claims are verified or clearly labeled as unverified, assumptions are explicit, and the deliverable does not smuggle speculation into the category of fact. *Verification status: Not verified. Human review required.* is not a disclaimer that lives on page 27; it is a posture that shapes the entire artifact.

Second, good output is **traceable to inputs**. Traceability means that the firm can point from any material claim back to either (i) a fact provided by the client or an approved internal source, or (ii) an assumption with an owner and rationale, or (iii) an interpretation that is explicitly framed as judgment. Traceability is the antidote to "AI said so." It also prevents the common failure mode where a deck contains a number that everyone has seen before but no one can source.

Third, good output is **reviewable with clear owners**. Reviewability means the artifact is structured so that reviewers can efficiently identify what must be checked and what can be accepted as a draft. Clear ownership means that the organization can answer: who verified this claim, who reviewed this narrative, who approved this deliverable, and who is accountable for the template that generated this structure? Reviewability is not optional. It is what makes scale possible without reducing quality to the lowest common denominator.

Fourth, good output is **retained with evidence**. Consulting firms already retain workpapers in various forms, but AI introduces new failure modes: prompt-driven drift, hidden transformations, and reuse-based propagation. Good Level 5 output is retained alongside the artifacts that explain it: inputs (appropriately redacted), prompts or templates (identified by ID), model/workflow lineage, reviewer notes, approvals, and evidence attachments. Retention is not merely for litigation defense. It is for learning, quality improvement, and credible internal supervision.

Fifth, good output is **reproducible enough to explain to a client and to internal risk**. Reproducible does not mean deterministic in the mathematical sense. It means operational reproducibility: the firm can recreate the deliverable context and explain how it was formed. If a client asks why an assumption was used, the firm can show when it was introduced, who owned it, and what verification plan existed. If internal risk asks how a template changed, the firm can show the release notes and evaluation results. Good output therefore supports both external credibility and internal governance.

In a Level 5 firm, an engagement deliverable is not a standalone artifact. It is the visible layer of a deeper workpaper structure. The firm builds the habit of thinking: "the deck is the narrative; the workpapers are the proof of professionalism." AI does not change that principle. It intensifies it, because AI can increase the distance between appearance and substance.

This definition also clarifies the boundary between consulting craftsmanship and consulting governance. Craftsmanship is still required: the ability to structure an argument, to tell a coherent

story, and to communicate tradeoffs. Governance ensures that craftsmanship does not become manipulation and that speed does not become carelessness. The purpose of Level 5 is not to eliminate human skill. It is to preserve the value of human skill under AI acceleration.

**Artifact (Save This)**

**Minimum deliverable at Level 5: the Consulting Workpaper Bundle (CWB).** For any AI-assisted engagement artifact that influences client-facing conclusions, the firm retains a CWB containing:

a) **Facts provided:** client-provided facts and approved internal sources, referenced and time-stamped. At Level 5, facts are treated as controlled objects. Each fact is tagged with provenance (client interview, client document, internal dataset, approved research), date, and owner. "Common knowledge" is not a provenance class for consulting deliverables. If a fact matters, it is either sourced or explicitly treated as an assumption. This prevents the quiet drift where an early draft includes an unsourced claim that becomes "true" by repetition.

b) **Assumptions register:** each assumption has an owner, rationale, materiality, and test/verification plan. The register is not a list of vague uncertainties. It is an operational tool. Owners are accountable for verifying or revising the assumption. Materiality indicates which assumptions would change the recommendation or the scenario ranking. The test plan specifies what evidence would confirm or falsify it and by when. The register is used to block downstream stages when critical assumptions remain unexamined.

c) **Open questions:** what remains unknown and what would change the conclusion. Open questions are tracked as blockers, not as curiosities. Each question is linked to a decision dependency: which part of the deliverable would change if the answer changes. This prevents teams from producing "complete" decks that are complete only because they silently guessed the missing answers.

d) **Model/workflow lineage:** model ID, prompt template ID, workflow version, parameter settings. Level 5 assumes that model behavior can change over time and that templates can drift. Lineage makes it possible to reconstruct outputs, compare versions, and run regression tests. It also prevents teams from using unapproved model configurations for high-risk deliverables.

e) **Outputs + reviewer notes:** draft outputs, edits, approvals, and sign-offs. The CWB retains not only what the model produced, but what humans changed and why. Reviewer notes document issues discovered (unsupported claims, unclear tradeoffs, overconfident tone) and the resolution. Approvals indicate that the deliverable passed required gates before external delivery.

f) **Verification evidence:** calculations, citations to internal sources, meeting notes, and validation checks. Evidence is attached to claims. Quantitative claims include calculation workpapers. Qualitative claims include interview notes or documented reasoning. If evidence is unavailable, the claim is either removed, softened into an assumption, or marked *Verification status: Not verified. Human review required.* with an explicit plan to verify.

g) **Risk log:** flags (privacy, hallucination, scope creep, decision laundering) and mitigation steps. The risk log is the firm's memory of what nearly went wrong. It records when sensitive data was present, when hallucination risk was elevated, when the team was tempted to overstate certainty, or when scope creep threatened to turn analysis into unreviewed advice. Mitigations are recorded so that governance improves over time rather than resetting each engagement.

The CWB is the practical centerpiece of the Level 5 mental model because it makes governance tangible. It turns abstract principles into concrete artifacts that teams can produce, reviewers can inspect, and the firm can retain. It also sets a professional expectation: in a Level 5 organization, "we delivered a deck" is not a complete statement. The complete statement is: "we delivered a deck, and we retained the evidence bundle that makes it defensible."

Finally, the Level 5 mental model should be understood as a continuous discipline rather than a destination. The firm does not become Level 5 once and remain there automatically. New teams join. New templates are created. New client contexts introduce new sensitivities. Model behavior shifts. The operating system must therefore include feedback loops: evaluations that run when assets change, monitoring that detects drift, and incident response that drives improvement. Level 5 is sustained not by enthusiasm for AI, but by institutional habits that treat governance as part of the craft of consulting delivery.

## 5.3 What Level 5 CAN do and CAN'T do in consulting-company practice

### 5.3.1 What it CAN do (when governance is real)

---

**Risk & Control Notes**

**Capability.** Standardize the consulting lifecycle into a supervised pipeline: intake structuring; completeness checks; scenario comparison; draft deliverables; approval routing; recordkeeping packages; and continuous monitoring of quality and drift.

**Primary risks.** Systemic error scales quickly: template drift, unverified claims replicated across teams, inconsistent approvals, missing workpapers, hidden conflicts, and brittle "AI says so" rationales.

**Minimum controls.** Define stage gates; enforce structured outputs; require attestations and approvals; maintain an AI asset registry; run evaluation harnesses; monitor exceptions; retain CWBs for reconstruction and internal QA.

---

At Level 5, the phrase "AI can help" stops being a statement about convenience and becomes a statement about the firm's delivery system. The core capability is not that the firm can draft faster. The core capability is that the firm can *standardize* how it delivers work under uncertainty while preserving professional judgment, evidentiary discipline, and supervisory accountability. When governance is real, Level 5 can make the consulting lifecycle behave like a controlled production pipeline: the same categories of inputs are captured, the same quality checks are performed, the same kinds of risks are flagged, and the same recordkeeping artifacts are retained across teams and offices.

The first major capability is **standardized intake and scoping**. In consulting, scope drift is not merely a commercial issue; it is a governance risk. As the scope expands, teams often import

new data, new assumptions, and new implicit claims without updating the verification posture. Level 5 can implement a firm-wide intake object that forces scoping decisions to be explicit: what decision is being supported; what deliverable is being produced; what data classes are involved; what tools are approved for that class; what constitutes "in scope" facts; and what cannot be claimed without validation. This standardization does not reduce consultant autonomy; it reduces ambiguity. It ensures that every team begins with the same governance posture rather than retroactively discovering constraints after a draft has already propagated.

The second major capability is **completeness checks that are operational rather than rhetorical**. Many consulting failures begin as omissions: a missing constraint, an unstated assumption, a dependency that was never surfaced. Level 5 can enforce completeness checks by requiring structured outputs at each stage. For example, the pipeline can require that any draft deliverable include (i) facts provided, (ii) assumptions introduced, (iii) open questions, (iv) claims ledger or evidence map, and (v) verification plan. If a consultant attempts to move forward without these fields, the system blocks the transition. This transforms "be careful" into a measurable process, and it reduces the probability that a polished narrative is mistaken for a complete analysis.

The third capability is **scenario comparison under an explicit assumption register**. Consulting work frequently involves comparing alternatives rather than proving a single truth. AI is useful here because it can help structure scenarios, enumerate tradeoffs, and generate consistent narrative frames. The risk is that scenario analysis becomes a stylistic exercise where assumptions are implicit and therefore untestable. Level 5 can make scenario comparison safer by enforcing an assumption register as a first-class artifact: every scenario is defined by its assumption set; material assumptions are flagged; and the scenario ranking is explicitly conditional on those assumptions. This produces a healthier consulting posture: not "this is the answer," but "given these assumptions, these are the tradeoffs; here is what would change the conclusion."

The fourth capability is **draft deliverables with embedded verification posture**. In a Level 5 firm, AI can draft slides, memos, workplans, and meeting briefs, but the drafting system does not pretend that drafts are final. It automatically tags outputs as *Verification status: Not verified. Human review required.* and requires that any quantitative or benchmark-like statement be either linked to evidence or rewritten into an assumption. This matters because consulting deliverables often mix hard facts with plausible approximations. Without a system, teams can slide from approximation to assertion. With Level 5 governance, the default is conservative: anything not verified is explicitly framed as provisional, and the verification plan becomes part of the deliverable lifecycle.

The fifth capability is **approval routing and separation of duties that scale**. Traditional consulting relies on hierarchical review: manager review, principal review, partner sign-off. AI does not remove the need for review; it increases it, because the volume of draft material increases and the persuasive quality of drafts can mask errors. Level 5 can standardize review roles and required approvals by deliverable class. For example, an internal working draft may require peer review, while a client-facing deck may require engagement leader review and partner approval, with

a compliance or risk review if sensitive data or regulated claims are involved. Crucially, Level 5 makes approvals evidence-based: reviewers are not approving aesthetics; they are approving that the claims ledger is acceptable, that assumptions are explicit, that verification evidence exists for material facts, and that prohibited phrasing has been avoided.

The sixth capability is **recordkeeping packages as a default output of doing the work**. Consulting firms already maintain workpapers, but the discipline varies. AI-assisted delivery increases the need for consistency because the transformations are faster and more numerous. Level 5 can define the Consulting Workpaper Bundle (CWB) as the minimum recordkeeping unit and can automate its creation: run manifest (model ID, parameters, workflow version), prompts log (with redactions and hashes), structured outputs, reviewer notes, approval records, evidence attachments, and risk flags. This turns recordkeeping from a compliance burden into an operational backbone. It also enables learning: when issues occur, the firm can diagnose causes and improve templates and controls.

The seventh capability is **continuous monitoring of quality and drift**. At Level 5, the firm does not assume that yesterday's template is safe today. Monitoring is therefore part of the delivery system. The firm can track exceptions (missing evidence, repeated risk flags, frequent rewrites for overconfidence, repeated use of unapproved assets) and can identify drift patterns (templates that gradually become more assertive, outputs that increasingly omit open questions, teams that bypass stage gates). Monitoring is not surveillance for its own sake; it is a reliability mechanism. It detects where the system is failing quietly before it fails publicly.

The eighth capability is **safe reuse of consulting IP**. Consulting firms win by reuse: standard workplans, diagnostic frameworks, playbooks, and slide structures. AI can amplify that advantage by accelerating adaptation. The risk is that reuse becomes uncontrolled replication, where outdated or incorrect content spreads faster than the firm can correct it. Level 5 can support reuse safely by treating reusable assets as governed products: versioned templates with owners; evaluation harnesses that test templates against known failure modes; release notes for changes; and deprecation policies for retiring assets. In other words, Level 5 allows the firm to scale leverage without scaling error.

The ninth capability is **a disciplined boundary between drafting support and professional judgment**. Level 5 governance is not merely about tools; it is about preventing the outsourcing of responsibility. The firm can use AI to structure options, articulate arguments, and generate drafts, but it must still own the judgment that selects a path and the accountability that stands behind it. When governance is real, the system supports this boundary by design: outputs are framed as drafts; verification is mandatory; and approvals are required for reliance-bearing conclusions. This is not a limitation. It is the mechanism that preserves the firm's credibility and the integrity of its work.

Taken together, these capabilities mean that Level 5 is not a promise of superhuman consulting. It is a promise of *professional consistency*: the firm can deliver faster without becoming sloppy, and it can reuse assets without becoming brittle. It is the maturity level at which AI becomes a sustainable part of delivery rather than a temporary productivity spike.

### 5.3.2 What it CAN'T do (do not outsource responsibility)

Level 5 is powerful precisely because it creates an operating model, but that same power creates a dangerous temptation: to treat the existence of controls as a substitute for judgment. This section makes the boundaries explicit. A Level 5 consulting organization can be well-governed and still be wrong; it can be audit-ready and still make a poor strategic recommendation. Governance improves defensibility; it does not guarantee correctness. Therefore, Level 5 cannot be used as a justification for outsourcing responsibility to the system.

First, Level 5 **cannot replace engagement partner sign-off**. In consulting, the partner's signature is the firm's credibility. It is the point where responsibility is concentrated, not diffused. Level 5 can make partner review more efficient by structuring the evidence and highlighting risk flags, but it cannot remove the need for a qualified human to own the final deliverable. Any attempt to treat AI as the final reviewer is a category error: the system can generate, structure, and flag, but it cannot own the professional consequences of what is delivered.

Second, Level 5 **cannot serve as the final arbiter of facts**. Even with secure tools and strong templates, an AI system cannot independently verify that a number is correct, that a claim is current, or that a benchmark is applicable to a particular client context. Facts require verification against approved sources. Level 5 can enforce the requirement that facts be sourced, but it cannot eliminate the human obligation to confirm sources and interpret them correctly. A sourced claim can still be misleading if context is wrong. Governance does not replace judgment; it frames judgment so it can be reviewed.

Third, Level 5 **cannot create auditability by writing logs after the fact**. A common failure mode in rushed environments is retrospective documentation: teams produce a deck and then try to reconstruct the process later. Level 5 rejects this. Auditability is produced by the workflow itself: stage gates, structured outputs, approvals, and retention bundles created as the work progresses. If the firm allows "we will document later" norms, it is not operating at Level 5, regardless of what tools it has.

Fourth, Level 5 **cannot justify a recommendation by style**. AI can produce consulting-shaped language: crisp headlines, tight storylines, and confident executive prose. This can be mistaken for rigor. At Level 5, style is treated as a risk amplifier: the more persuasive the narrative, the greater the obligation to show evidence and to surface assumptions. A recommendation is justified by reasoning and evidence, not by a compelling slide title. The firm must resist the substitution of narrative polish for analytical substance.

Fifth, Level 5 **cannot bypass client confidentiality**. Consulting involves sensitive data: strategy, financials, operations, organizational issues. No operating model can justify violating confidentiality norms or contractual obligations. Level 5 therefore cannot allow teams to ingest sensitive client data into unapproved tools, to share prompts externally, or to retain client data outside controlled retention systems. Even "anonymized" data can be re-identifying if context is rich. The model may be useful, but the firm's confidentiality obligations are non-negotiable.

Sixth, Level 5 **cannot generate external benchmarking as if verified**. Consulting deliv-

erables often include market sizes, peer comparisons, cost benchmarks, and "industry norms." AI can generate plausible benchmarks quickly, which is exactly why this is dangerous. Unless the firm has validated sources and an approved process for using them, AI-generated benchmarks must be treated as placeholders or hypotheses, not as facts. Level 5 can help by forcing a claims ledger and evidence mapping, but it cannot conjure verified benchmarks from thin air.

Seventh, Level 5 **cannot eliminate bias, persuasion risk, or value judgments**. Consulting recommendations involve normative choices: tradeoffs between growth and risk, efficiency and resilience, centralization and autonomy. AI can help articulate tradeoffs, but it cannot decide values. Moreover, AI-generated language can subtly bias a decision by framing options asymmetrically. Level 5 governance must therefore require explicit alternatives and balanced framing, but it cannot guarantee neutrality. Humans must review not only for correctness but for framing integrity.

Eighth, Level 5 **cannot reduce all risk to procedures**. Procedures are necessary, but they can be gamed. A team can fill in a template and still be sloppy. A reviewer can sign off perfunctorily. A system can produce logs that no one reads. Therefore Level 5 cannot be treated as a mechanical compliance exercise. The operating model must be paired with culture: an expectation that evidence matters, that assumptions are surfaced, that review is real, and that shortcuts are unacceptable for client-facing conclusions.

Ninth, Level 5 **cannot prevent all incidents**. Even mature systems will produce occasional failures: a missed assumption, an outdated benchmark, a confidentiality near-miss, an overconfident phrasing. The purpose of Level 5 is not to claim immunity. It is to ensure that incidents are detected quickly, contained, corrected, and learned from, with clear accountability and improvement of controls. A firm that believes it is incident-proof is, by definition, not governed.

These "cannot" statements are not a pessimistic add-on; they are the definition of professional seriousness. Consulting firms that treat Level 5 as a pathway to autonomy will eventually face the same outcome: institutionalized error, reputational damage, and loss of trust. Consulting firms that treat Level 5 as a pathway to disciplined scale can use AI to increase productivity while strengthening the evidentiary and supervisory backbone of delivery.

---

**Risk & Control Notes**

**Hard stop examples (Level 5).**

1. Auto-sending client deliverables (or sharing drafts externally) without human verification and documented approval.
2. Treating AI-generated benchmarks, market sizes, costs, or "industry norms" as factual without validated sources.
3. Using AI outputs to rationalize a pre-chosen recommendation (decision laundering) instead of documenting reasoning.
4. Ingesting sensitive client data into unapproved tools or retaining it outside firm recordkeeping controls.

---

The practical function of the hard stops is to make governance enforceable. They turn vague

caution into explicit prohibitions that can be embedded into workflow design: external delivery requires approval artifacts; benchmarks require evidence mapping; scenario framing requires alternatives; sensitive data requires approved tools and retention. Level 5 is not a set of nice ideas. It is a system that prevents predictable failures by refusing to proceed when the firm cannot demonstrate defensibility.

## 5.4 Core workflow patterns for Level 5 (organization-level discipline)

Level 5 is where governance stops being a set of principles and becomes a delivery system. The firm is no longer asking teams to "use good judgment" and hoping that good judgment survives time pressure, staffing changes, and template reuse. Instead, the firm designs workflows so that disciplined behavior is the default and undisciplined behavior is visible, difficult, and costly. In practice, that means the firm adopts a small number of core workflow patterns that recur across engagements and service lines, each with embedded controls, stage gates, and recordkeeping artifacts. These patterns are not meant to be clever. They are meant to be enforceable.

A useful way to interpret the patterns below is to think of them as the consulting equivalent of a well-run factory line, with one decisive difference: consulting is not manufacturing. The inputs are ambiguous, the outputs are probabilistic, and judgment is inseparable from delivery. Therefore the goal is not to eliminate variation; it is to bound variation and to make the sources of variation explicit. Level 5 patterns do exactly that. They force the firm to say, in operational terms, what counts as an acceptable input, what counts as an acceptable output, what must be verified before delivery, and who must sign off.

### 5.4.1 Pattern A: Standardized intake → scoping gate (minimum-necessary data)

Level 5 begins before any drafting, analysis, or slide building occurs. It begins with intake. Intake is often treated as administrative overhead, but in an AI-assisted environment intake is a primary risk control. The firm is about to place client context into a system that can transform it rapidly into plausible narratives. If the intake is sloppy, the outputs will be persuasive and wrong, or persuasive and unsafe, or both. If the intake is disciplined, the workflow can enforce sensible boundaries.

A Level 5 intake object is a structured schema, not an email thread. It is the firm's canonical representation of what this engagement is and is not. At minimum, it captures: (i) the engagement objective (decision to support, not just topic), (ii) the deliverable type (memo, deck narrative, workplan, workshop brief), (iii) the client context and constraints, (iv) the sensitivity classification of data (low/medium/high, or whatever scheme the firm uses), (v) the permitted tool list for that classification, (vi) the approved sources that may be used (internal knowledge bases, client documents, approved research repositories), and (vii) explicit "not in scope" boundaries.

The minimum-necessary principle is central here. Teams often over-share context because it feels helpful. But in governed consulting, "helpful" must be balanced against confidentiality and necessity.

Level 5 intake enforces that discipline by requiring teams to separate what is essential for drafting from what is merely convenient. For example, it should be possible to draft a workplan and an interview guide without embedding sensitive financials. It should be possible to structure hypotheses and issue trees without pasting proprietary process documents. When the system forces teams to state what they need and why, it reduces the likelihood that sensitive material is unnecessarily introduced into the workflow.

The scoping gate is the hard stop that follows intake. It is the firm's first formal checkpoint, and it has one job: prevent work from starting when the engagement has not been defined well enough to govern. The scoping gate is passed only if the intake is complete enough to support safe drafting. Typical required conditions include:

a) the objective and deliverable type are explicitly defined;
b) sensitivity classification is selected and consistent with the inputs provided;
c) the permitted tool list is acknowledged (and unapproved tools are prohibited);
d) key facts provided are enumerated as facts (not embedded in narrative);
e) known missing facts are enumerated as open questions;
f) "not in scope" boundaries are recorded;
g) reviewer and approver roles are assigned for the deliverable class.

This gate matters because it creates a discipline that consulting teams otherwise struggle to maintain: early clarity. Without a scoping gate, teams tend to draft first and clarify later. In an AI environment, drafting first becomes especially dangerous because the first draft can quickly become the anchor for the entire engagement narrative, even if it was generated from partial or misclassified inputs. The Level 5 scoping gate prevents that anchoring by blocking progression until the engagement is governable.

### 5.4.2 Pattern B: Drafting with embedded verification plans (facts ≠ assumptions)

Once the engagement passes intake and scoping, the firm enters drafting. Level 5 drafting is not the same as "write a deck." It is the disciplined generation of intermediate artifacts that are *Verification status: Not verified. Human review required.* by default and that include a built-in plan to move from draft to defensible deliverable. This pattern is the backbone of governance-first consulting: it ensures that every output is accompanied by a verification posture and an evidence plan.

The heart of the pattern is a strict output schema that appears in every stage output. The exact schema may vary by firm, but the core fields must remain stable: facts provided; assumptions; open questions; draft output; verification plan; evidence needed; risk flags; and handoff notes. The firm should treat this schema as non-negotiable. It is the mechanism that prevents persuasive error. If the model produces a fluent paragraph, the workflow requires the paragraph to be linked to facts or to be recast as an assumption. If the model produces a number, the workflow requires a source or forces it into the "unverified placeholder" category with an explicit plan to verify.

Embedded verification plans also require ownership. A verification plan that has no owner is a

polite fiction. Therefore, Level 5 drafting requires that each verification item be assigned to a role or individual (even if initially as a placeholder). Ownership makes verification operational rather than aspirational. It also supports stage gating: if the verification owner has not provided evidence, the deliverable cannot proceed past the gate.

The principle "no evidence, no claims" is most important for quantitative statements: benchmarks, market sizes, cost baselines, and comparative performance claims. Level 5 drafting treats these as high-risk by default. The workflow can permit teams to use placeholders, but only under explicit labeling. For example, a draft may include "[Placeholder: verify TAM from approved source]" rather than a fabricated number. The system can also require a "claims ledger" view that enumerates each material claim and records whether evidence exists. This makes review more efficient and makes it harder for unsupported claims to slip into polished output.

A subtle but crucial dimension of this pattern is language discipline. Drafting systems should enforce conservative phrasing when facts are unknown. They should avoid words that imply certainty ("will," "proves," "guarantees," "best," "always") unless the claim is genuinely verified and appropriate. They should encourage conditional phrasing ("may," "could," "under these assumptions," "indicative," "requires confirmation") to preserve the correct epistemic posture. In Level 5, this is not a stylistic preference; it is a governance control. Overconfident language increases client misunderstanding and increases the probability that internal reviewers accept the content without scrutiny.

Finally, embedded verification plans convert uncertainty into a work program. Consulting work inevitably begins with unknowns. The difference between disciplined consulting and sloppy consulting is not whether unknowns exist; it is whether unknowns are documented, owned, and resolved in a controlled manner. Pattern B makes that distinction operational.

### 5.4.3   Pattern C: Separation of duties (producer vs reviewer vs approver)

Level 5 assumes that a single operator is a single point of failure, especially when AI can generate large volumes of plausible content quickly. Pattern C therefore formalizes two-person integrity. The firm defines distinct roles, assigns them explicitly, and designs workflows so that one person cannot unilaterally generate and ship client-facing content.

At minimum, the pattern separates three roles:

a) **Producer:** runs the workflow, drafts outputs, and assembles the workpaper bundle.
b) **Reviewer:** evaluates reasoning posture, evidence mapping, assumption discipline, and compliance with the output schema; identifies unsupported claims and required revisions.
c) **Approver:** typically the engagement leader or partner; confirms that the deliverable is fit to share externally and that required approvals and evidence are present.

These roles may be played by different people depending on engagement size, but the separation must be real. The reviewer must be independent enough to challenge content. The approver must have authority and accountability. If the same individual plays all roles, the firm may still generate

logs, but it will not achieve meaningful governance.

This pattern exists primarily to combat two failure modes.

The first is **single-operator drift**. When one person controls both generation and approval, there is a natural tendency to accept outputs that align with their intuition, especially under time pressure. AI can magnify this by producing persuasive drafts that feel "good enough." A reviewer creates friction in the correct place: where persuasion might outrun verification.

The second is **tool charm**. AI tools can be charismatic. They produce confident text, crisp bullet points, and consulting-style storylines. This can create a psychological shortcut: "it sounds right." Separation of duties forces a second set of eyes to evaluate whether it *is* right, or at least whether it is appropriately framed as unverified.

Operationally, Pattern C is enforced by stage gates and attestations. The producer can draft, but cannot mark the artifact as ready for client delivery. The reviewer can approve internal quality, but cannot bypass partner sign-off for external delivery. The approver cannot sign without the reviewer's notes and the presence of required evidence. These constraints turn a cultural norm into a workflow reality.

In addition, Pattern C supports training and capability development. If reviewers regularly flag the same types of issues (unsupported benchmarks, missing assumptions, overconfident tone), the firm can update templates and training materials. Governance therefore becomes a feedback loop rather than a policing function.

### 5.4.4   Pattern D: Reuse safely: templates + playbooks + evals + change control

Consulting firms run on reuse. Reuse is leverage, and leverage is profitability. AI increases the ability to reuse because it can adapt content quickly to new contexts. But reuse is also how errors become systemic. Pattern D is therefore the Level 5 pattern that treats reusable AI assets as controlled products.

The first element is **standardized prompt templates and expected output schemas**. A Level 5 firm does not let teams invent new prompt structures for high-risk deliverables ad hoc. It provides approved templates that produce outputs in required schemas. This standardization reduces variability in quality and makes review more consistent. It also makes evaluation feasible: you cannot test what you cannot standardize.

The second element is **test cases**. For every reusable asset (prompt template, workflow, playbook), the firm maintains a small suite of test cases that represent common engagement scenarios and known failure modes. These are not primarily about model "performance" in the abstract. They are about governance outcomes: does the template separate facts and assumptions; does it avoid prohibited claims; does it generate open questions; does it produce an evidence map; does it maintain conservative tone; does it avoid leaking sensitive data; does it fail safely when inputs are missing?

The third element is an **evaluation harness and regression suite**. The firm must be able to run the test cases whenever an asset changes: when a template is modified, when a model version

changes, when the workflow is updated, or when a new use case is added. Regression testing is what prevents template drift from quietly undermining governance. Without regression tests, the firm will gradually accumulate small changes that each seem harmless but collectively erode discipline.

The fourth element is **change control**. Change control means: no template is updated and pushed firm-wide without an owner, a documented rationale, evaluation results, and release notes. It also means that assets have versions and that teams can identify which version they used on an engagement. If a client issue arises, the firm can trace it to a specific asset version and can decide whether to roll back, patch, or retire the asset. This is essential for accountability and for organizational learning.

The fifth element is **deprecation**. Firms tend to add assets but rarely remove them. At Level 5, deprecation is a governance requirement: outdated templates and playbooks are retired, blocked, and removed from default catalogs. Otherwise, the firm's "approved" library becomes a museum of inconsistent practices, and teams will pick whichever asset is convenient rather than whichever is governed.

Pattern D is the mechanism that enables safe scale. It allows the firm to distribute AI-enabled workflows widely while preserving a baseline of defensibility. It is also the mechanism that prevents the organization from becoming brittle. A brittle firm is one that cannot change because it fears unintended consequences. A Level 5 firm can change because it has tests, owners, and rollback paths.

### 5.4.5   Pattern E: Exceptions, escalations, and kill-switch logic

Even the best-designed operating model will face exceptions. Level 5 governance is not credible unless it includes a plan for exceptions that is as operational as the plan for normal delivery. Pattern E defines how the firm stops, escalates, and, when necessary, disables assets or workflows.

The first step is **exception classification**. The firm should define a small set of exception categories that matter operationally, such as:

a) **Confidentiality exception:** sensitive client data was entered into an unapproved tool, or retention was uncontrolled.

b) **Factual integrity exception:** suspected hallucination or unsourced quantitative claim made it into a draft deliverable.

c) **Decision process exception:** evidence of decision laundering, suppressed alternatives, or biased framing.

d) **Process exception:** missing approvals, missing workpaper bundles, bypassed stage gates.

e) **Client dispute exception:** client challenges a claim, interpretation, or representation, triggering reconstruction needs.

Classification matters because it determines escalation paths. Not every exception requires the same response. A minor tone issue can be resolved by a rewrite. A confidentiality issue may require immediate containment, internal reporting, and client communication posture.

The second step is **stop rules**. Stop rules define when work must halt. This is where Level 5 becomes real. Examples include: sensitive data detected in an unapproved context; missing evidence for a material quantitative claim; deliverable class requiring partner approval but approval missing; or repeated template failures in evaluation. Stop rules prevent "we will fix it later" behavior, which is one of the most common routes to serious incidents.

The third step is **escalation**. Escalation is not merely telling a partner. It is routing the exception to the right control owner: risk/compliance for confidentiality; knowledge management for template drift; engagement leadership for process exceptions; and firm leadership for systemic issues. Escalation is also time-bound. A system that escalates but does not resolve is not governed; it is bureaucratic.

The fourth step is **kill-switch logic**. A Level 5 firm must be able to disable an asset quickly when it is unsafe. This can mean disabling a specific template, blocking a workflow for a deliverable class, or restricting a model configuration while investigation occurs. Kill-switch logic is a hallmark of mature systems because it acknowledges reality: occasionally the correct action is to stop using something immediately, not to debate.

The fifth step is **incident response and learning**. When an exception occurs, the firm should be able to reconstruct what happened using the Consulting Workpaper Bundle: inputs, lineage, outputs, reviewer notes, approvals, and evidence. The point is not blame; it is improvement. The firm updates templates, improves evaluation tests, strengthens stop rules, and retrains teams. Over time, this turns exceptions into a source of organizational learning rather than a source of repeating embarrassment.

Pattern E also protects client relationships. When a client challenges a claim, the firm that can reconstruct its process quickly and transparently will retain credibility even if a correction is needed. The firm that cannot reconstruct will appear careless, regardless of whether the content was "mostly right." In consulting, credibility is often more fragile than the content itself.

Taken together, Patterns A through E define the operational heart of Level 5. They turn governance from a set of principles into a lived delivery system: intake that controls what enters, drafting that preserves epistemic discipline, separation of duties that prevents single-point failure, reuse governed by evaluation and change control, and exception handling that stops problems before they scale. A firm that implements these patterns is not merely using AI. It is governing AI as part of consulting professionalism.

## 5.5 Mini-cases (Level 5): organization-level delivery with audit-ready evidence

The mini-cases in this book are intentionally repetitive. They are not intended to impress the reader with variety; they are intended to force maturity. In Levels 1 through 4, the same scenarios teach progressively stronger disciplines: drafting discipline at Level 1, explicit reasoning at Level 2, stage-gated multi-step workflows at Level 3, reusable assets and evaluation at Level 4. At Level 5,

repetition serves a different purpose. It tests whether the firm can run the *same operating model* across different engagement types and still produce work that is reconstructable, reviewable, and defensible under pressure.

This section therefore treats each mini-case as a stress test for the organization. The cases are not about finding the perfect strategic answer. They are about proving that the consulting company can deliver an answer responsibly: with clear scoping, controlled inputs, structured outputs, explicit assumptions, verification plans, evidence bundles, approvals, retention, and post-delivery learning. The most important output of Level 5 is not the memo or the deck. It is the **Consulting Workpaper Bundle (CWB)** that makes the memo or the deck defensible.

Across the four cases, the firm should demonstrate the same principles in different contexts:

1. **Intake discipline:** minimum-necessary data; sensitivity classification; tool permissions; in-scope and out-of-scope boundaries.
2. **Epistemic discipline:** facts are not assumptions; unverified claims are labeled; uncertainty is visible.
3. **Evidence discipline:** material claims map to sources, calculations, or explicit assumptions with owners and tests.
4. **Accountability discipline:** separation of duties; reviewer notes; approvals are recorded.
5. **Operational discipline:** templates are versioned; evaluation harnesses exist; exceptions are escalated; retention is consistent.

When these disciplines are present, AI becomes a controlled accelerator. When they are absent, AI becomes a persuasion engine that can scale error.

### 5.5.1   Case 1: Market entry strategy (industry expansion)

**Scenario frame.** A mid-sized industrial services company has a strong position in its home market and is considering expansion into a new geography and adjacent service line. The board is divided: some directors want growth, others want risk reduction and capital discipline. The CEO wants a recommendation in eight weeks, with a decision memo and a supporting deck for the board meeting. Stakeholder alignment is critical because the expansion would require a new operating footprint, new channel partnerships, and potentially acquisitions.

The engagement begins with uncertainty. The client provides some facts (current revenue, margin ranges, high-level customer segmentation, high-level competitive landscape, internal capability assessments) but lacks reliable market sizing and has limited customer research in the target geography. The client also has confidentiality constraints: only a subset of financials can be shared, and some documents are restricted to the engagement team under strict handling rules.

A Level 5 firm treats this scenario as a classic environment for AI-enabled drafting and structuring *with strong governance.* The risk is high because market entry deliverables are often heavy on externally flavored claims: market sizes, growth rates, competitor moves, "industry norms," and success factors. AI can generate those claims easily, which is precisely why they must be governed.

The firm must prevent the deck from becoming an elegant fiction.

The intake object for this case should therefore include: the decision criteria (risk-adjusted growth, time-to-scale, strategic fit, capability gaps, capital requirements, execution risk); the timeline and key decision points; the permitted internal sources (client documents, interview notes, approved research repositories); prohibited sources (ad hoc web browsing unless validated through approved research process); the sensitivity classification; and the explicit rule that any market sizing, competitor claims, or benchmarks must be treated as *Verification status: Not verified. Human review required.* unless sourced to approved evidence.

**Level 5 outputs the firm can produce (with governance).**

a) **A stage-gated workplan with explicit decision points and evidence requirements.** The workplan is not simply a timeline. It is a governance artifact. It specifies stages (intake and scoping; hypothesis and issue tree; evidence gathering; scenario design; synthesis; review and approval; board pack finalization), and it ties each stage to required artifacts and gates. For example, the firm may define a "market evidence gate" that must be passed before any market-size numbers can appear in client-facing materials. The workplan also assigns owners and reviewers, and it defines what constitutes acceptable evidence (approved research notes, validated interview summaries, documented calculations). This is where Level 5 turns ambiguity into controlled execution.

b) **A scenario comparison memo with assumptions register and verification plan (no fabricated market data).** The scenario memo is produced early and iterated. It is explicitly conditional. It describes 2–4 entry scenarios (e.g., greenfield expansion, partnership-led entry, acquisition-led entry, or a phased pilot) and frames tradeoffs. The memo separates facts provided, assumptions, and open questions. Critically, it includes an assumption register with owners and verification plans: for example, assumptions about achievable customer acquisition rates, channel access, regulatory constraints, and operational ramp-up. Any placeholders for market sizing are labeled and tied to the evidence plan. The goal is to avoid the common failure mode where a persuasive narrative appears before evidence exists.

c) **A CWB package linking each claim to a source or to an assumption with owner/test.** The CWB is the deliverable that proves professionalism. It includes a claims ledger for the board pack: each major claim (market attractiveness, competitive intensity, capability gaps, investment requirements, expected ramp) is labeled as fact, assumption, or interpretation; facts are tied to sources; assumptions have owners and test plans; interpretations are tied to structured reasoning and documented tradeoffs. Reviewer notes record what was challenged and how it was resolved. Approvals indicate that the partner signed off on the final board pack under a defined verification posture.

In a Level 5 firm, the board deck is therefore the visible output, but the CWB is the operational truth. It is what allows the firm to answer the most damaging client question: "Where did you get that?" If the firm cannot answer that question crisply, the deliverable is not defensible, no matter how beautiful the slides are.

---

**Artifact (Save This)**

**Minimum deliverable (Case 1): Market Entry CWB checklist.**

☐ Intake schema completed (objective, deliverable, sensitivity class, permitted sources/tools, out-of-scope).

☐ Issue tree and hypothesis list retained with version history.

☐ Assumptions register with owners, materiality flags, and verification plans.

☐ Open questions list tied to decision dependencies and deadlines.

☐ Claims ledger for the board pack mapping each claim to evidence or assumption.

☐ Evidence map with attached calculations, interview notes, and approved research references.

☐ Reviewer notes and partner approval record for external delivery.

☐ Final deck/memo retained with lineage (template IDs, workflow version, timestamps).

---

### 5.5.2 Case 2: Cost transformation (SG&A + operations)

**Scenario frame.** A diversified services company faces margin compression and has committed to a cost transformation program. The CEO wants a 12–18 month plan that delivers measurable savings while maintaining service quality. The CFO demands a credible baseline and an auditable savings model. HR is concerned about workforce implications and messaging. The organization is politically sensitive: cost programs are prone to skepticism and internal resistance, and claims of savings can become reputational liabilities if they are not realized.

Cost transformations are a classic consulting domain where AI can both help and harm. AI can accelerate the drafting of workplans, diagnostics, and communications. It can also produce confident but incorrect savings claims, and it can flatten nuanced tradeoffs into generic recommendations. At Level 5, the firm's job is to prevent the engagement from turning into a narrative machine. The work must remain calculation-driven, assumption-transparent, and auditable.

The intake object here must explicitly capture: the baseline definition (which costs, which period, which allocation rules), data sources (general ledger extracts, payroll, procurement data), confidentiality constraints, who can see what, and the governance policy that no savings claim can appear in client-facing materials without a calculation workpaper and a signed-off assumptions register. It must also flag HR/legal sensitivities: workforce-related content requires minimum-necessary exposure and may require additional approvals depending on the firm's policy and client's constraints.

**Level 5 outputs the firm can produce (with governance).**

a) **A "claims ledger" tying cost-savings statements to calculations, sources, and sign-offs.** The claims ledger is the central control artifact. Every savings claim is decomposed into a formula, with inputs sourced to baseline data and assumptions explicitly stated. For example, a procurement savings claim references spend categories, current unit costs, negotiated discounts, implementation timelines, and adoption curves. Each element is tied to an evidence source or an assumption with an owner. The ledger records reviewer sign-off for each claim category.

This prevents the most common cost transformation failure: promising savings that are not supported by traceable math.

b) **A deliverable template with prohibited phrasing (avoid guarantees) and required disclaimers.** Cost programs carry persuasion risk. A consulting deck can easily imply certainty about savings realization, timeline, and organizational ease. Level 5 templates therefore include prohibited phrasing rules: avoid "will deliver" unless backed by evidence and caveats; avoid "guarantee" language; avoid presenting point estimates without ranges where uncertainty is material; avoid implying that workforce impacts are trivial. The template also includes required disclaimers and verification posture language: savings are estimates contingent on execution, client decisions, and adoption. This is not legal hedging; it is truth-in-consulting.

c) **A QA checklist for spreadsheet integrity, assumptions, and rounding/aggregation risk.** Cost programs live in spreadsheets, and spreadsheets are where errors become expensive. A Level 5 firm maintains a QA checklist that reviewers actually use: reconcile totals to baseline; check time periods and inflation assumptions; verify headcount and wage-rate sources; check double-counting across levers; review rounding and aggregation; verify that scenario ranges are coherent; and confirm that implementation costs are included rather than silently excluded. The checklist is retained in the CWB as evidence that the firm did not merely generate narrative.

A mature Level 5 implementation also integrates approvals: certain categories of claims (workforce reductions, facility closures, vendor changes) may require additional client-side stakeholder review before being placed in a board pack. The firm's operating model should anticipate that and build it into stage gates rather than leaving it to last-minute negotiation.

---

**Risk & Control Notes**

**Case 2: verification gates.** Cost transformations require hard verification gates because the content is inherently consequential and often politically charged. At minimum:

1. **Baseline gate:** no savings claim proceeds without a reconciled baseline definition and source documents.

2. **Model gate:** no lever-level savings claim proceeds without a calculation workpaper and assumption owner.

3. **Workforce sensitivity gate:** workforce-related content requires minimum-necessary exposure and appropriate approvals.

4. **Timeline feasibility gate:** implementation timing assumptions must be reviewed for realism; point estimates should be treated conservatively when uncertain.

These gates exist to prevent persuasive but fragile savings narratives from becoming client commitments.

### 5.5.3 Case 3: Capital allocation (divest, invest, return capital)

**Scenario frame.** A publicly visible company is facing pressure from investors to improve returns. The board is evaluating whether to divest a non-core business, invest aggressively in a growth segment, or return capital through buybacks and dividends. Each option has reputational and governance implications. The CEO needs a decision-ready board book with clear alternatives, tradeoffs, and a defensible rationale. The environment is high-stakes: leaks can move markets, and poor framing can trigger investor backlash.

This case is particularly important at Level 5 because it reveals the limits of AI and the importance of governance posture. Capital allocation is not merely analysis; it is governance. It involves fiduciary duties, board dynamics, stakeholder communications, and sometimes regulatory sensitivity. Even when the consulting firm is not providing legal advice, it is operating in a context where process legitimacy matters as much as analytical content.

The intake object must therefore be unusually strict. It should classify data as highly sensitive, specify tool permissions tightly, define retention rules, and include explicit boundaries: the firm will not treat AI-generated external claims as verified; the firm will not produce any client-facing statement implying certainty about market reactions; and the firm will not allow the system to "recommend" an action without human-owned reasoning.

**Level 5 outputs the firm can produce (with governance).**

a) **A decision book structure (no recommendation by AI): alternatives, tradeoffs, assumptions.** The key Level 5 output here is structural. The firm produces a decision book template that forces balanced framing: each alternative is presented with benefits, costs, risks, execution requirements, and failure modes. Assumptions are explicit and linked to evidence or to a verification plan. The language is careful to distinguish analysis from advocacy. AI can assist by drafting consistent sections and ensuring completeness, but the firm prohibits "AI recommendation" outputs. The deliverable is a disciplined decision-support document, not an automated judgment.

b) **A risk register covering model risk, data risk, and narrative risk.** In capital allocation work, narrative risk is unusually high: how options are framed can influence board decisions and external perception. The Level 5 risk register therefore includes: model risk (hallucination risk, drift, template issues), data risk (confidentiality, provenance, inconsistencies), and narrative risk (bias, asymmetric framing, overconfidence). It also includes mitigation steps: mandatory alternative framing checks, conservative language rules, explicit uncertainty disclosures, and strict approvals.

c) **A sign-off routing pack (who approves what, with what evidence).** This pack formalizes accountability. It defines which sections require which approvals: financial model assumptions reviewed by finance workstream lead; narrative framing reviewed by engagement leader; risk disclosures reviewed by the firm's risk/compliance function (if applicable to the consulting company's governance posture); partner sign-off before any board-facing document is shared. The routing pack is retained as part of the CWB, making it clear that the firm treated the

deliverable as high-stakes and supervised it accordingly.

The distinctive Level 5 lesson in Case 3 is that governance is not an add-on for "regulated industries" only. Capital allocation work inside any major company is governed work because it affects stakeholders and can have material consequences. Therefore, the consulting company's operating model must scale up controls when stakes increase: stricter tool permissions, stricter retention rules, tighter separation of duties, and more explicit sign-offs.

### 5.5.4 Case 4: Operating model redesign (org + process + tech)

**Scenario frame.** A large organization is redesigning its operating model after years of ad hoc growth. The client wants clarity on roles, decision rights, processes, and enabling technology. The engagement involves workshops, interviews, process mapping, and the creation of new RACI matrices and governance forums. Change management is central: the client is not only selecting a new design but also persuading leaders to adopt it.

Operating model redesign is where consulting firms often rely heavily on reusable assets: RACI templates, governance models, process maps, and workshop materials. AI can accelerate these dramatically by drafting variations and summarizing workshop outputs. The risk is that the firm ships "generic operating model" content that is not anchored in the client's context, or that it fails to retain evidence of how decisions were reached, which becomes problematic when stakeholders disagree.

At Level 5, this case is an opportunity to demonstrate safe reuse and institutional learning. The firm should treat its operating model playbook as a governed asset with evaluation harnesses and version history. It should also treat workshop and interview outputs as sensitive data requiring controlled handling and retention.

**Level 5 outputs the firm can produce (with governance).**

a) **A controlled "playbook" asset (templates + checklists + test cases) with version history.** The firm's operating model playbook is a prime candidate for Level 5 governance. It includes standardized templates (RACI, governance forum design, decision-rights matrices), checklists (what must be clarified before proposing design changes), and test cases (common failure modes, such as ambiguous accountability or overloaded decision forums). The playbook is versioned, owned, evaluated, and released under change control. AI can be used to adapt templates to client contexts, but only through governed workflows that preserve the required output schema and verification posture.

b) **A monitoring plan: drift detection in templates, compliance phrasing, and missing workpapers.** Operating model work tends to produce many artifacts: workshop notes, drafts, revised RACIs, and stakeholder communications. Drift is a real risk: templates can gradually lose required disclaimers, or teams can omit evidence bundles under schedule pressure. A Level 5 firm therefore defines monitoring signals: missing CWB components; repeated bypass of stage gates; template usage without evaluation; frequent changes to governance language; and repeated

stakeholder disputes that suggest unclear reasoning. Monitoring allows the firm to catch process degradation early.

c) **A training and certification plan for staff use (minimum competency + prohibited uses).** Operating model work is often staffed with mixed experience levels. Level 5 therefore includes a training plan that makes governance enforceable: staff must demonstrate they can use the structured output schema; they can distinguish facts from assumptions; they can produce and retain CWBs; they understand confidentiality rules and tool permissions; and they know the hard stops (e.g., do not paste sensitive workshop transcripts into unapproved tools). Certification is not about bureaucracy; it is how the firm ensures that governance does not depend on a few experts.

The Level 5 lesson in Case 4 is that reuse is not merely an efficiency tactic; it is a governance risk surface. The more the firm reuses assets, the more it must treat those assets like products: owned, tested, monitored, and controlled. When done well, Level 5 governance allows the firm to reuse structure and discipline while still tailoring content to the client's reality.

Across all four cases, the pattern is consistent. Level 5 does not promise that the firm will always produce the best strategic answer. It promises that the firm will produce an answer responsibly: with explicit uncertainty, traceable claims, documented assumptions, enforced review, and reconstructable evidence. That is what "organization-level delivery with audit-ready evidence" means in a consulting company. The content may change from case to case, but the operating model does not.

## 5.6 Risks and controls taxonomy for Level 5 (minimum operating model)

Level 5 is where the consulting firm stops thinking about AI risk as a collection of individual mistakes and starts treating it as a systems problem. At small scale, errors are painful but containable: a consultant writes an overconfident email; a manager misses an unsupported claim; a slide includes an unsourced number that gets corrected before delivery. At firm scale, the same errors can become systemic. They propagate through reuse, accelerate through templates, and become embedded into "how we do things here." That is why Level 5 governance is not optional. It is the operating model that prevents the organization from industrializing error.

This taxonomy is intentionally practical. It does not attempt to be academically complete. Its purpose is to define the minimum set of risks a consulting company must expect at scale and the minimum set of controls required to remain defensible. The guiding assumption is conservative: *Verification status: Not verified. Human review required.* by default, human ownership of judgment, and auditable evidence as a first-class deliverable. The taxonomy also reflects the key Level 5 insight: controls must be *embedded in the workflow.* Policies that live in PDFs do not scale. Controls that block or shape work do.

### 5.6.1 Risk categories (what can go wrong at firm scale)

1. **Confidentiality and client data governance:** improper tool use; uncontrolled retention; cross-client contamination.

   Confidentiality risk is the most obvious, but at Level 5 it becomes more complex than "someone pasted sensitive data into the wrong place." Consulting firms operate across multiple clients, multiple engagements, and multiple internal repositories. AI increases the probability that information moves in ways humans cannot easily track. The firm faces several sub-risks:

   *Improper tool use.* Teams may use unapproved tools for convenience, especially when under deadline pressure. The risk is not only that data leaves the firm; it is also that the firm cannot demonstrate where data went, under what controls, and with what retention posture.

   *Uncontrolled retention.* AI workflows generate artifacts: prompts, intermediate drafts, structured outputs, logs, and evaluation results. Without defined retention rules, the firm may retain too much (creating confidentiality exposure) or retain too little (destroying defensibility and reconstruction capability). "We do not store anything" is rarely credible at Level 5, because the firm must retain enough evidence to supervise and reconstruct.

   *Cross-client contamination.* The most dangerous confidentiality failure at firm scale is not a single leak; it is the accidental mixing of client-specific knowledge across engagements. This can occur when consultants reuse templates that contain embedded client examples, when internal repositories are not segmented properly, or when AI-assisted drafts inadvertently incorporate prior engagement language. Even if no explicit confidential fact is copied, the appearance of cross-client contamination undermines trust.

   Confidentiality risk is therefore governed by more than privacy training. It requires tool permissions, data classification, minimum-necessary input rules, workspace segregation, and enforceable workflows that prevent mixing.

2. **Systemic factual error:** plausible but wrong statements propagated through templates and reuse.

   AI makes it easy to produce "consulting-shaped facts": market sizes, benchmarks, maturity scores, industry norms, cost baselines, and peer comparisons. At small scale, a wrong number might be caught by a subject matter expert. At firm scale, wrong statements can become standardized if they are embedded into reusable assets. The firm then faces an existential reputational risk: it is not merely wrong once; it is wrong repeatedly and consistently.

   Systemic factual error typically emerges through three mechanisms:

   *Template drift.* A template begins with placeholders and cautious language, but over time small edits remove qualifiers, add specificity, or introduce example numbers that were meant to be illustrative. The template becomes more assertive without becoming more verified.

   *Copy-forward normalization.* Consultants copy slides from prior decks. AI accelerates this by generating variants that feel "consistent." Over time, the firm forgets which claims were verified and which were inherited.

   *Source ambiguity.* Teams cannot distinguish between claims sourced from approved research and

claims generated as plausible filler. If provenance is not explicit, everything becomes equally "true" by the time it reaches the final deck.

The result is a firm-scale failure mode: the organization produces credible-looking nonsense at industrial speed.

3. **Decision laundering:** using AI to legitimize weak reasoning or pre-decided outcomes.

Decision laundering is the governance risk that most directly threatens consulting integrity. It occurs when AI is used to create the appearance of rigorous analysis while the underlying decision was pre-decided, politically driven, or insufficiently examined. AI makes decision laundering easier because it can generate structured narratives, issue trees, and tradeoff language that feel analytical even when they are not anchored.

At Level 5, decision laundering is not merely an ethical issue; it is an operating model risk. If the firm's pipeline rewards polished output over explicit reasoning and evidence, then decision laundering becomes a rational behavior under deadline pressure. Teams will use AI to produce plausible justification rather than to expose uncertainty and alternatives.

Decision laundering often appears as:

a) suppressed alternatives (only one path is framed as "reasonable");

b) asymmetric framing (the preferred option gets benefits and minimal risks; others get risks and minimal benefits);

c) retrofitted assumptions (assumptions are chosen to support the conclusion rather than tested);

d) lack of falsification plan (no explicit "what would change our mind").

If this risk is not controlled, Level 5 turns into "institutionalized persuasion" rather than disciplined consulting.

4. **Process breakdown:** missing approvals; inconsistent stage gates; inability to reconstruct workpapers.

Level 5 requires a stage-gated delivery system. The risk at scale is that the firm has a designed process but does not have an enforced process. Process breakdown occurs when teams bypass gates, approvals become rubber stamps, or workpapers are incomplete. This is often caused by organizational pressure: deadlines, partner expectations, understaffing, or client demands for speed.

The most damaging consequence is the inability to reconstruct. When a client challenges a claim, when internal risk asks for evidence, or when a quality incident occurs, the firm cannot answer the basic questions:

a) What exactly was produced?

b) Which inputs were used?

c) Which assets and versions were involved?

d) Who reviewed and approved it?

e) What evidence supported the claims?

Without reconstruction capability, the firm appears careless even if the content was mostly correct. Consulting credibility is partially procedural: clients trust firms that can show their work.

5. **Model/asset drift:** changes in model behavior, prompt templates, or inputs degrade quality over time.

Drift is inevitable. Models change, even if the firm does not request changes. Prompt templates evolve. Consultants adapt patterns. Client inputs change. Over time, these shifts can degrade output quality or governance posture. A template that used to produce cautious language might begin producing overconfident phrasing after a model update. A prompt that used to elicit explicit assumptions might begin producing implicit assertions. A workflow that used to flag missing evidence might begin missing it because the input format changed.

Drift risk is distinct from factual error. A system can drift into:

a) lower completeness (fewer open questions surfaced),
b) higher assertiveness (stronger claims with less evidence),
c) weaker separation of facts and assumptions,
d) increased sensitivity leakage (more client details included by default),
e) inconsistent tone across teams (making review harder).

Without monitoring and regression testing, drift will be invisible until it becomes a client incident.

6. **IP and provenance risk:** mixing client IP, firm IP, and public content without clear provenance. Consulting companies have multiple overlapping IP domains: client confidential materials, firm proprietary frameworks, public-domain concepts, and vendor-provided content. AI can blur boundaries. It can rewrite client materials in a way that retains substance; it can generate language that resembles proprietary frameworks; it can incorporate public content without clear citation posture. At Level 5, the risk is not only legal. It is reputational and operational: if provenance is unclear, the firm cannot confidently reuse assets, cannot confidently publish thought leadership, and cannot confidently separate client deliverables from firm IP.

Provenance risk includes:

a) accidental inclusion of client-specific content in reusable templates;
b) accidental reuse of firm proprietary language in contexts where it should not appear;
c) untracked incorporation of public content that creates plagiarism or attribution issues;
d) ambiguity about what was generated vs what was sourced.

In a governed consulting organization, provenance must be explicit enough to support reuse decisions, publication decisions, and confidentiality obligations.

### 5.6.2 Controls (minimum viable Level 5 control set)

The control set below is intentionally framed as minimum viable. It is not the full suite a large firm may eventually adopt, but it is the lowest threshold at which Level 5 can credibly claim to be an

operating model rather than a collection of intentions. These controls must be both documented and operationalized: they must exist in policy *and* in workflows, tools, and stage gates.

1. **AI asset registry:** owners, purpose, permissible use, sensitivity class, version, deprecation date.

   The AI asset registry is the foundation of organizational governance. If the firm does not know what AI assets exist, it cannot govern them. The registry treats AI-related artifacts as production assets: prompt templates, workflow definitions, evaluation harnesses, model configurations, reusable playbooks, and any automation wrappers that structure inputs/outputs. A minimum registry entry includes:

   a) asset ID and name (stable identifier),
   b) owner (accountable role or person),
   c) purpose and intended deliverable classes,
   d) permissible use statement (what it can do; what it must not do),
   e) sensitivity class allowed (what data can be used),
   f) required output schema (facts/assumptions/open questions, etc.),
   g) version history and release date,
   h) evaluation status (last test date; test suite used; pass/fail summary),
   i) deprecation date or review cadence (to prevent stale assets).

   The registry is not a catalog for curiosity. It is the enforcement hook for governance. If an asset is not in the registry, it is not approved for use in client-impacting work.

2. **Stage gates and approvals:** defined checkpoints; role-based sign-offs; escalation paths.

   Stage gates are where governance becomes real. They are the points in the workflow where progress is blocked unless required artifacts and approvals exist. At Level 5, stage gates are not optional reminders in a slide footnote; they are enforced transitions in the delivery pipeline. A minimum stage-gated model includes:

   a) **Intake/scoping gate:** sensitivity classification; permitted tools; objectives; out-of-scope boundaries; owners assigned.
   b) **Claims readiness gate:** quantitative claims either sourced or labeled as assumptions with test plans; claims ledger created.
   c) **Review gate:** independent reviewer notes recorded; key issues resolved; risk flags addressed.
   d) **Approval gate:** engagement leader/partner sign-off recorded before any external sharing.
   e) **Retention gate:** CWB completed and stored per policy immediately after delivery.

   Escalation paths are part of the gate design: confidentiality exceptions route to risk/compliance; systemic template issues route to knowledge management; repeated drift indicators route to the asset owner and governance committee.

3. **Standard output schemas:** facts/assumptions/open questions + verification plan + evidence map.

   Output schemas are the simplest and most effective control for preventing persuasive error.

They force the system to keep epistemic categories separate. In a Level 5 firm, client-impacting outputs must follow a standard schema that includes:

a) facts provided (with provenance tags and timestamps),

b) assumptions (with owners and materiality),

c) open questions (with decision dependencies),

d) draft output (clearly labeled *Verification status: Not verified. Human review required.* until verified),

e) verification plan (who will verify what, by when),

f) evidence map (claim-to-source mapping),

g) risk flags (privacy, hallucination risk, decision laundering risk).

Schemas make review scalable. They also make monitoring feasible because the firm can measure missing fields, repeated risk flags, and deviations from required structures.

4. **Evaluation harness:** test cases for each reusable asset; regression tests before release.
Evaluation harnesses prevent systemic error from spreading. A harness is not a one-time demo; it is a repeatable set of tests that can be executed whenever an asset changes. For each reusable asset, the firm maintains a small test suite that checks governance outcomes: separation of facts/assumptions, presence of open questions, conservative language, absence of prohibited claims, and correct handling of missing inputs.
Regression tests are run:

a) before deploying a new template version,

b) after a model update or configuration change,

c) when adding a new deliverable class or new use case,

d) periodically (scheduled) to detect drift over time.

Evaluation results are stored and referenced in the asset registry. If an asset fails, it is not released. This is how Level 5 prevents template drift from quietly eroding governance.

5. **Change control:** template/model updates require review, test results, and release notes.
Change control is the bridge between innovation and safety. Without it, firms either freeze (fear of breaking things) or drift (constant untested change). At Level 5, change control means:

a) every change has an owner and a documented rationale;

b) evaluation harness results are attached to the change request;

c) reviewer approval is required for release;

d) release notes record what changed and what risks were considered;

e) rollback plans exist for critical assets.

Change control also includes deprecation: the firm retires assets that are outdated, unsafe, or redundant. This prevents the approved library from becoming a confusing sprawl that undermines consistency.

6. **Recordkeeping and retention:** CWBs retained per policy; ability to reconstruct "what was produced, when, and why".

Recordkeeping is the mechanism that turns governance into defensibility. The CWB is the minimum unit of retention for AI-assisted consulting work that influences client-facing conclusions. Retention policies must specify where CWBs live, who can access them, and how long they are kept. The firm should be able to reconstruct:

a) the inputs and their provenance (with appropriate redaction),

b) the workflow and asset versions used,

c) the draft outputs and edits,

d) reviewer notes and approvals,

e) the evidence supporting material claims,

f) the risk flags and mitigations.

This capability is not merely defensive. It supports learning, quality improvement, and credible supervision.

7. **Monitoring and exceptions:** dashboards for missing artifacts, high-risk outputs, and drift indicators.

Monitoring is the control that keeps Level 5 alive over time. Without monitoring, governance decays. A minimum monitoring set includes dashboards or periodic reports that flag:

a) missing CWB components (e.g., no claims ledger, no evidence map),

b) high-risk outputs (e.g., repeated quantitative placeholders, repeated decision laundering flags),

c) drift indicators (e.g., reduced open questions, increased assertive language),

d) unapproved asset usage (assets not in registry),

e) repeated bypass of stage gates or approvals.

Exceptions trigger escalation and, when necessary, kill-switch actions: disabling an asset, blocking a workflow path, or requiring additional approvals until the issue is resolved. Monitoring is therefore not passive observation; it is an operational feedback loop.

A key Level 5 principle is that these controls reinforce each other. The asset registry enables change control. Change control requires evaluation harness results. Evaluation harnesses depend on standardized schemas. Schemas make monitoring possible. Monitoring identifies exceptions, which feed back into asset updates and training. Recordkeeping makes it possible to reconstruct incidents and improve controls. The system is designed to be coherent rather than fragmented.

---

**Artifact (Save This)**

**Minimum Standard for Safe Use at Level 5 (printable operating checklist).**

☐ Each AI asset in use appears in the registry with an owner and permitted-use statement.

☐ Every engagement run produces a CWB with lineage, reviewer notes, approvals, and evidence.

☐ Quantitative claims are either sourced and verified or explicitly labeled as assumptions with a test plan.

☐ Stage gates block delivery when required facts/evidence or approvals are missing.

☐ Reusable templates/playbooks have an evaluation harness and release notes for every change.

☐ Exceptions (privacy, hallucination risk, client disputes) trigger escalation and (if needed) kill-switch.

---

If a consulting company can meet this minimum standard consistently, it has crossed the threshold into Level 5 behavior. It can credibly say that AI is governed as part of delivery, not treated as an informal productivity tool. The firm may still refine and expand controls, but it has the essential operating model: assets are known and owned, workflows are gated, outputs are structured, changes are tested, evidence is retained, and exceptions are managed. That is what it means, in practice, to be an AI-capable consulting organization without sacrificing defensibility.

## 5.7 Prompt patterns and exercises (copy/paste)

At Level 5, prompt patterns are not "tips." They are controlled interfaces into the firm's delivery operating system. The primary objective is not to get a clever answer; it is to produce auditable artifacts that preserve governance posture under deadline pressure. Therefore, every Level 5 prompt pattern has three properties: (i) it forces a structured output schema that separates facts from assumptions, (ii) it embeds verification planning and evidence needs directly into the output, and (iii) it is designed to be retained inside the Consulting Workpaper Bundle (CWB) with clear lineage (template ID, version, and reviewer notes).

A second Level 5 principle is that prompts should be written as if they will be reviewed by a skeptical reader later: a partner, an internal risk reviewer, or a client. This is not paranoia. It is professionalism. If the prompt and its output cannot survive a retrospective review, they should not be used for client-impacting work. This is also why the prompts below are deliberately conservative. They explicitly forbid recommendations, require explicit uncertainty, and impose a *Verification status: Not verified. Human review required.* posture by default.

Finally, these prompts are intended to be copy/paste-ready but not context-free. The consulting company should version them, store them in its AI asset registry, test them via an evaluation harness, and apply change control when modifying them. A prompt template is an organizational asset. Treat it like one.

### 5.7.1 Prompt Template 1: Organization-level intake and scoping (Level 5)

This template is the default entry point for Level 5 work. It is designed to produce a scoping artifact that can pass the first stage gate: the team can proceed to drafting only if the intake is complete enough, sensitivity classification is correct, and "not in scope" boundaries are explicit. The output is structured JSON so it can be stored, diffed, and reviewed.

**Prompt (Copy/Paste)**

```
ROLE: You are a governance-first drafting assistant inside a consulting firm.
You do NOT make recommendations. You produce structured drafts and workpaper artifacts.

OUTPUT FORMAT (STRICT JSON):
{
  "facts_provided": [],
  "assumptions": [],
  "open_questions": [],
  "analysis": [],
  "draft_output": [],
  "verification_plan": [],
  "evidence_needed": [],
  "risk_flags": [],
  "handoff_notes": [],
  "verification_status": "Not verified"
}


ENGAGEMENT INTAKE (MINIMUM-NECESSARY):
- Client industry: [...]
- Objective (draft-only): [...]
- Deliverable type: [memo / deck narrative / meeting brief / workplan]
- Sensitivity class: [low / medium / high]
- Facts provided (ONLY these are facts): [...]
- Constraints: [time, tone, prohibited claims, internal policy notes]
- Required stage gates: [scoping, evidence, review, approval]

TASK:
1) Produce a scoping draft that restates objective, non-goals, required inputs, and verification gate
2) Produce a first-pass workplan (stages + owners placeholders + evidence per stage).
3) Produce the initial assumptions register and open questions list.
```

**How to use this template well (and safely).**

1. **Keep inputs minimum-necessary.** Do not paste sensitive client documents. Summarize at the level required for scoping and drafting.
2. **Treat the scoping output as a gate artifact.** The team should not proceed to drafting until

the scoping draft is reviewed and missing facts are identified.

3. **Force specificity on "not in scope."** Many consulting failures are scope failures disguised as analysis. The prompt should explicitly list what the deliverable will not do.

4. **Assign owners quickly.** The output may include placeholder owners. Replace them with real roles before progressing.

### 5.7.2 Prompt Template 2: Claims ledger (tie each claim to evidence)

At Level 5, the claims ledger is one of the highest-leverage controls. It converts a narrative into a reviewable object: each claim is classified, mapped to evidence, assigned an owner, and rewritten conservatively if unverified. This template is designed to be run on any draft memo or deck narrative before it reaches partner approval.

---

**Prompt (Copy/Paste)**

```
Take the draft below and produce a "claims ledger" table in JSON:
- claim_text
- claim_type: [fact / assumption / interpretation]
- evidence_required
- evidence_present: [yes/no]
- verification_owner
- risk_if_wrong
- rewrite_if_unverified (more conservative wording)

Rules:
- If a claim contains a number, benchmark, market size, cost estimate, or "industry norm",
  assume evidence_present="no" unless the draft explicitly cites an approved source.
- If the claim is not verifiable (forward-looking), classify it as assumption or interpretation,
  and require an assumptions register entry with an owner and test plan.
- Do NOT invent sources. Do NOT add external facts. Rewrite conservatively when unverified.

Draft:
[PASTE DRAFT]
```

---

**What reviewers should look for.**

1. Claims that look factual but are actually assumptions (e.g., "customers will adopt," "implementation will take six months").

2. Claims that embed benchmarks without provenance (e.g., "best-in-class," "top quartile," "industry standard").

3. Language that implies certainty where evidence is missing (rewrite into conditional form).

4. Any claim that would embarrass the firm if challenged (those require immediate verification or removal).

### 5.7.3 Exercises (20–30 minutes each)

These exercises are designed to train Level 5 behavior: not prompting skill in isolation, but disciplined production of auditable artifacts. Each exercise should produce a small bundle of outputs that can be reviewed: a claims ledger, a stage-gated workflow, a mini evaluation harness, or a policy note. The objective is to build organizational muscle memory.

1. **Claims ledger conversion and conservative rewriting.** Pick a past (anonymized) engagement memo or deck narrative. Run Prompt Template 2 to generate a claims ledger. Then take every claim where `evidence_present="no"` and rewrite it into an explicitly-assumptive form that includes: (i) the assumption statement, (ii) what evidence would be needed to verify it, and (iii) what would change if the assumption is wrong. Deliverables: the claims ledger JSON; a revised draft with conservative language; and a short reviewer note summarizing the top three risk-if-wrong items.

2. **Design a stage-gated workflow with three mandatory stops.** For the market-entry mini-case, design a workflow with three mandatory hard stops:

   a) **Missing facts stop:** the workflow cannot proceed until the intake object includes the decision criteria and the list of open questions that would change the conclusion.

   b) **Evidence stop:** the workflow cannot proceed to client-facing narrative until quantitative claims are mapped in a claims ledger and either verified or rewritten as assumptions.

   c) **Partner approval stop:** the workflow cannot be externally shared until reviewer notes and partner sign-off are recorded.

   Deliverables: a one-page workflow diagram or bullet workflow; gate criteria; and the list of required CWB artifacts at each gate.

3. **Create a minimal evaluation harness for a reusable template.** Build a five-test evaluation harness for a reusable "market entry memo" template. For each test, define:

   a) the test prompt (with minimal dummy facts),

   b) the expected output schema fields that must appear,

   c) pass/fail criteria in plain English (e.g., "must produce at least 5 open questions," "must not include numbers without marking unverified," "must separate assumptions explicitly"),

   d) a common failure mode the test is designed to detect.

   Deliverables: a test suite document; pass/fail rubric; and a note on how the firm would run regression tests before releasing a template update.

4. **Write a one-page incident posture: "What we do when AI is confidently wrong."** Draft a short internal policy note that defines the firm's posture when an AI-assisted artifact contains a confident error. The note must include:

   a) **Detection:** how errors are identified (review, monitoring flags, client challenge).

   b) **Containment:** immediate actions (stop sharing, freeze the asset version, preserve evidence).

   c) **Escalation:** who is notified (engagement leader, risk/compliance, knowledge management).

d) **Correction:** how the deliverable is corrected and re-approved (claims ledger update, revised evidence map).

e) **Client communication posture:** conservative language, transparency, and remediation without over-admission.

f) **Learning:** how the evaluation harness and templates are updated to prevent recurrence.

Deliverable: the one-page policy note suitable for inclusion in the firm's governance playbook and retention as part of a CWB when an incident occurs.

These prompt patterns and exercises are intentionally designed to feel procedural. That is the point. Level 5 is the maturity level where process is the artifact. A consulting company that can run these patterns repeatedly, retain the workpapers, and improve templates through evaluation and change control is not merely using AI—it is governing AI as part of delivery.

## 5.8 Conclusion: operating Level 5 and maintaining discipline over time

Level 5 is not a finish line. It is an operating posture that must be maintained. The organization can implement an AI asset registry, define stage gates, publish templates, and launch dashboards—and still drift back into informal behavior as soon as deadlines tighten and attention shifts. This is the central paradox of Level 5: the very success of the system can make it feel routine, and routine systems are where governance quietly decays. When the firm gets used to producing structured outputs and workpaper bundles, it becomes tempting to treat them as boxes to check rather than as evidence of real discipline.

This conclusion therefore emphasizes a practical reality: Level 5 is a culture of defensibility enforced by an operating model. The firm's goal is not to maximize the amount of AI in the workflow. The firm's goal is to maximize the reliability of delivery while preserving professional judgment. When Level 5 is functioning, the firm can move fast without becoming reckless, reuse assets without spreading error, and scale delivery without scaling harm. When Level 5 is not functioning, the firm can still produce impressive-looking deliverables, but it cannot prove how they were produced, why they are reasonable, or what is verified versus assumed. That is not maturity. That is risk disguised as competence.

### 5.8.1 Summary of main takeaways

The core takeaway of Level 5 is simple and worth repeating: **Level 5 is a delivery operating system plus evidence.** It is not a tool selection. It is not a platform license. It is not a prompt library. It is the integration of policy, workflow, controls, and recordkeeping into the act of doing consulting work.

First, Level 5 reframes AI from a drafting helper into a controlled component inside a supervised pipeline. The system begins with standardized intake: minimum-necessary data, sensitivity

classification, and explicit boundaries. This is where governance becomes real because it controls what enters the workflow and under what conditions.

Second, Level 5 makes epistemic discipline enforceable. Facts are separated from assumptions and open questions in every structured output. Anything unverified is labeled *Verification status: Not verified. Human review required.* by default. Quantitative claims are treated as high-risk by default: if there is no evidence, there is no claim. This is not a philosophical statement; it is a stage-gate rule that blocks progression.

Third, Level 5 scales accountability through separation of duties. A single operator cannot generate and ship client-facing content. Producers, reviewers, and approvers have defined roles; approvals are recorded; and the partner remains responsible for external-facing deliverables. This prevents the organization from outsourcing responsibility to "the system" and prevents the psychologically powerful failure mode of tool charm.

Fourth, Level 5 enables safe reuse by treating reusable assets as governed products. Prompt templates, workflows, and playbooks have owners, versions, permitted-use statements, and deprecation paths. They are evaluated through test suites and regression harnesses before release. Change control exists not to slow innovation but to prevent template drift from quietly eroding governance posture.

Fifth, Level 5 treats recordkeeping as a primary deliverable. The Consulting Workpaper Bundle is not a compliance afterthought; it is the mechanism that makes consulting defensible. The firm can reconstruct what was produced, when, and why. It can show which inputs were used, which assets and versions were involved, what assumptions were made, what evidence supported each material claim, what edits occurred, and who approved delivery. In a world where AI accelerates drafting and can produce persuasive text quickly, reconstruction capability is not optional. It is the firm's credibility insurance.

Sixth, Level 5 recognizes that drift is inevitable and therefore builds monitoring and exception handling into the operating model. The firm watches for missing artifacts, missing approvals, high-risk outputs, and drift indicators. Exceptions trigger escalation, and if needed, kill-switch logic disables unsafe assets quickly. This is not paranoia. It is how mature systems behave: they assume failures will occur occasionally and design for rapid detection and containment.

Finally, Level 5 clarifies the purpose of governance: it is not "more AI." It is **more defensibility**. The firm does not win by producing the most AI-generated text. The firm wins by producing client-impacting work that can withstand scrutiny—from clients, internal risk, and the firm's own professional standards. Governance is what allows reuse to scale without scaling harm.

### 5.8.2   What comes next (beyond Level 5)

Beyond Level 5, the organization does not move into autonomy. It moves into continuous improvement. The system becomes more complete, more tested, and more resilient, but it remains governed and human-owned. The firm's maturity is measured not by the absence of errors but by how quickly errors are detected, contained, corrected, and learned from.

In practice, "beyond Level 5" is a set of reinforcing loops:

**Evaluation coverage expands.** Early Level 5 implementations often test a subset of high-value templates and workflows. Over time, the firm expands evaluation harness coverage: more deliverable classes, more failure modes, more adversarial tests, and more scenario variation. The objective is not to achieve perfection; it is to reduce unknown unknowns by converting them into known tests. When a failure occurs, the firm does not merely correct the deliverable. It adds a test so the failure becomes less likely to recur.

**Assets evolve under disciplined change control.** The firm's asset library will grow, but growth must be curated. Mature firms implement a "golden path": a small set of default approved assets that cover the most common engagement needs. New assets are introduced cautiously, with owners and evaluation. Old assets are deprecated aggressively. The firm resists the natural tendency toward sprawl because sprawl undermines consistency and makes review harder.

**Monitoring becomes predictive rather than reactive.** Early monitoring often focuses on missing artifacts and obvious violations. Over time, the firm can monitor more subtle indicators of drift: reductions in open questions, rising rates of unverified quantitative placeholders, patterns of reviewer corrections, repeated bypass attempts, or increases in assertive language. Monitoring becomes a quality signal, not just a compliance alarm.

**Training and certification tighten.** Level 5 does not rely on informal learning. It operationalizes competency. Staff are trained not only to use templates but to understand the governance logic: why facts and assumptions must be separated; why claims require evidence; why stage gates exist; how to respond to exceptions. Certification should be lightweight but real: individuals demonstrate they can produce a compliant output schema, build a claims ledger, and assemble a CWB. As the firm learns, training is updated and integrated into onboarding and role progression.

**Incident response becomes a source of organizational learning.** Mature organizations treat incidents as input into improvement, not as embarrassing anomalies to be hidden. When an AI-assisted error occurs, the firm preserves evidence, reconstructs what happened, corrects the deliverable, communicates conservatively, and updates assets and tests. Over time, incident response becomes faster and more consistent, and the firm's credibility increases because it behaves like a serious professional organization rather than a tool experimenter.

**Controls mature without becoming burdensome.** A common fear is that governance will slow delivery. The goal is the opposite: embed controls so that disciplined work is faster than undisciplined work. Over time, the firm can reduce friction by improving tooling: automatic CWB assembly, integrated claims ledgers, standardized evidence mapping, and streamlined approval routing. The firm's competitive advantage becomes not "AI speed" but "safe speed."

The most important point is that the firm should resist the idea that the next step is autonomy. Consulting is judgment under uncertainty, and judgment is accountable. Level 5 is the maturity level where the organization proves that it can use AI as a disciplined accelerator while retaining human ownership. Beyond Level 5 is simply better Level 5: broader coverage, tighter feedback loops, and stronger learning discipline.

---

**Artifact (Save This)**

**Level 5 "stay here" checklist (operational maintenance).**

1. **AI asset registry is current; deprecated assets are retired and blocked.** Owners are assigned; permissible use statements are maintained; assets without owners are removed; and the firm prevents use of outdated templates by default. "Approved" must mean something operational.

2. **Evaluation harness runs on a schedule and on every change (release discipline).** Regression tests are executed for every template and workflow change, and they are also executed periodically to detect drift. Test results are stored and referenced in release notes. Failing assets are not released.

3. **Monitoring detects missing CWBs, missing approvals, and drift patterns.** Dashboards or periodic reports surface exceptions: missing workpaper bundles, bypassed stage gates, repeated high-risk outputs, and changes in output behavior that suggest drift. Monitoring triggers escalation and corrective action.

4. **Training and certification are enforced; exceptions are logged and reviewed.** Staff demonstrate minimum competency for using governed workflows. Exceptions are not ignored; they are logged, reviewed, and used to update training and controls. High-risk roles and deliverable classes have tighter requirements.

5. **The firm can reconstruct any external-facing claim: inputs, reasoning posture, evidence, edits, and approval.** Reconstruction is the ultimate test. When challenged, the firm can show what was known, what was assumed, what was verified, who reviewed it, and why it was reasonable to present. If reconstruction is not possible, the firm is not truly operating Level 5.

---

If there is a single sentence that captures Level 5 maintenance, it is this: **governance must be easier than improvisation.** The firm should design the operating model so that the fastest way to produce a client-ready deliverable is also the most defensible way. When that is true, discipline sustains itself. When it is not true, discipline erodes, and the organization returns to ad hoc behavior where AI amplifies speed but also amplifies risk.

A consulting company that can maintain Level 5 over time earns a distinctive advantage. It is not merely efficient; it is reliable. It can scale teams rapidly while preserving quality. It can reuse assets confidently without fear of hidden defects. It can survive scrutiny because it can show its work. In a market where clients increasingly care not just about answers but about how answers were produced, that reliability becomes a strategic asset. Level 5 is therefore not the end of the story. It is the point where the firm becomes capable of learning and improving without losing control.

# Bibliography

[1] National Institute of Standards and Technology (NIST). *Artificial Intelligence Risk Management Framework (AI RMF 1.0).* NIST AI 100-1, 2023.

[2] International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC). *ISO/IEC 23894:2023: Artificial intelligence — Guidance on risk management.* ISO/IEC, 2023.

[3] International Organization for Standardization (ISO). *ISO 31000:2018: Risk management — Guidelines.* ISO, 2018.

[4] International Organization for Standardization (ISO) and International Electrotechnical Commission (IEC). *ISO/IEC 27001:2022: Information security, cybersecurity and privacy protection — Information security management systems — Requirements.* ISO/IEC, 2022.

[5] Organisation for Economic Co-operation and Development (OECD). *Recommendation of the Council on Artificial Intelligence.* OECD/LEGAL/0449, 2019.

[6] European Parliament and Council of the European Union. *Regulation (EU) 2024/1689 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act).* Official Journal of the European Union, 2024.

[7] National Institute of Standards and Technology (NIST). *The NIST Cybersecurity Framework (CSF) 2.0.* NIST CSWP 29, 2024.

[8] Committee of Sponsoring Organizations of the Treadway Commission (COSO). *Internal Control — Integrated Framework.* COSO, 2013.

# Appendix A

# Notebook Index (Companion Colab Notebooks)

This appendix lists the companion notebooks for each chapter. Each notebook run generates an auditable bundle (run manifest, prompts log, risk log, deliverables).

**Repository path (GitHub):** `https://github.com/alexdibol/ai-consulting/tree/main/notebooks`

| Chapter | Notebook (file) and focus |
| --- | --- |
| Chapter 1 | `chapter_1.ipynb`: Level 1 drafting workflows, redaction hygiene, facts/assumptions/open-questions separation, basic logging. |
| Chapter 2 | `chapter_2.ipynb`: Level 2 structured reasoning, alternatives/tradeoffs, assumption register + tests, scenario frames, questions-to-verify. |
| Chapter 3 | `chapter_3.ipynb`: Level 3 multi-step workflow with gates + QA, review packet bundling, reconstructability artifacts. |
| Chapter 4 | `chapter_4.ipynb`: Level 4 reusable assets + evaluation harnesses, versioning, controlled release, monitoring + rollback. |
| Chapter 5 | `chapter_5.ipynb`: Level 5 operating model simulation: intake → routing → workflow → QA → approval → recordkeeping. |