# Multi-agent LLM System Blueprint for Telecoms

November 2025

**Silent Comet Team** 





**O döcomo Business** 

Part of the NaaS Accelerator Program



# **Table of contents**

1	Intro	oduction	(
	1.1	Why this framework exists	6
	1.2	Framework overview	6
	1.3	Reading guide	ç
2	Exe	cution lifecycle	1(
	2.1	Adapted workflow at a glance	10
	2.2	Phase 1: Scoping (Project definition)	11
	2.3	Phase 2: Data (Knowledge & context)	12
	2.4	Phase 3: Modeling (Agent design, ML integration & evaluation)	14
	2.5	Phase 4: Deployment (Production, ML Ops and monitoring)	16
	2.6	Putting it all together	18
3	Arcl	nitectural viewpoints	20
	3.1	Data & Model Lifecycle viewpoint	2
	3.2	Governance, Risk, and Compliance viewpoint	22
	3.3	Runtime/Process (Execution Topologies) viewpoint	23
	3.4	Capability/Development viewpoint	24
	3.5	Physical/Infrastructure viewpoint	25
	3.6	End-to-End Scenario viewpoint	26
	3.7	Applying viewpoints incrementally	27
4	-	metrics & KPIs	3
	4.1	Metrics by lifecycle phase	31
	4.2	Metrics by architectural viewpoint	32
	4.3	Telecom-specific KPIs	33
	4.4	Metric interdependencies and trade-offs	34
	4.5	Application guidance	36
5	Fran	nework in action	37
	5.1	Multi-agent reference architecture	37
	5.2	Use case validation	41
	5.3	Implementation roadmap	44
	5.4	Conclusions	45
6	Sum	nmary	47
	6.1	Framework validation: problems solved	47
	6.2	Final recommendations	50
	6.3	Closing perspective	5





Table of contents 3

A	appendices	52
Α	ML Workflow Blueprint - Key Aspects	52
	A.1 Phases descriptions	52
	A 2 Do's and don'ts by stage	59





# List of Figures

1.1	Architectural viewpoints populated by lifecycle phases	7
1.2	The measurable flywheel: each iteration compounds value and organizational learning	8
	Lean agent system lifecycle showing iterative workflow phases with feedback loops Three-layer architecture showing Foundation (data/models), Orchestration	10
	(agents/workflow), and Integration (APIs/interfaces)	12
5.1	Three-Layer Agent Architecture for Telecom	38





# **List of Tables**

3.1	Viewpoint adoption progression aligned with lifecycle phases	28
3.2	Viewpoint-to-phase integration matrix showing artifact flows	29
4.1	Lifecycle phase metrics summary	31
4.2	Viewpoint metrics summary	32
4.3	Network operations KPIs	33
	Customer operations KPIs	34
	Leading-to-lagging metric relationships	35
5.1	Framework adaptation across use cases	43
6.1	Execution validation: Delivery and integration obstacles addressed through framework design	47
6.2	Governance validation: Risk mitigation through architectural implementation	49





## 1 Introduction

## 1.1 Why this framework exists

Telecom operators face a predictable AI adoption trap: proof-of-concept demos that never reach production, escalating costs without measurable ROI, and compliance gaps that surface too late. The root cause isn't purely technical - it's structural. Teams lack a shared language between Technology, Operations, Risk, and Finance, which leads to diffuse scope, brittle prototypes, and governance treated as an afterthought.

The pattern repeats across operators: an executive sponsors an AI pilot for network troubleshooting or service order automation. The team builds an impressive demo in weeks. Then reality hits nobody defined baseline KPIs, the system retrieves outdated network documentation, production traffic reveals 30% hallucination rates, cost per interaction exceeds manual processing, and compliance demands audit trails that don't exist. Six months later, the pilot is quietly shelved, stakeholders are fatigued, and the organization concludes "AI isn't ready for telecom."

The underlying failure modes appear consistently: undefined ROI metrics, data overreach delaying time-to-value, fragile prototypes lacking governance, uncontrolled costs eroding business cases, and disconnected ML assets that can't integrate into operational workflows. That said, these obstacles are addressable through intentional design choices validated across production deployments.

This framework offers a pragmatic adoption route: start with one scoped, high-value process, prove measurable uplift quickly, keep the architecture intentionally lean and expand only when KPIs validate the investment.

## 1.2 Framework overview

The framework integrates two complementary perspectives that evolve together, aligning Technology, Operations, Risk/Compliance, and Finance around a shared language of value, control, and scalability:

Four-phase execution lifecycle (Chapter 2) drives project through Scoping  $\rightarrow$  Data  $\rightarrow$  Modeling  $\rightarrow$  Deployment, with feedback loops enabling continuous refinement.

**Six architectural viewpoints** (Chapter 3) capture decisions and govern evolution, transforming lifecycle artifacts into durable governance assets (see Figure 1.1).







These aren't parallel perspectives - they're integrated lenses on the same system. Lifecycle phases produce artifacts for architectural views; viewpoints define what is important to look at and document in each of the stages.



Figure 1.1: Architectural viewpoints populated by lifecycle phases

**Lifecycle phases** produce concrete artifacts: Scoping establishes KPI baselines and boundaries; Data curates GOLD knowledge sets with quality gates; Modeling builds evaluated prompts, tools, and ML integrations; Deployment monitors quality, cost, and governance with rollback capability.

Viewpoints provide the interpretive and organizational framework for how these artifacts are





used, combined, and understood. They transform these artifacts into governance assets. Teams may start with three core viewpoints (Data & Model Lifecycle, Runtime/Process, End-to-End Scenario) and layer in others as the system matures. Every production trace maps back to versioned decisions across viewpoints, which helps prevent architectural drift.

This integration creates a measurable flywheel where each iteration builds on the last, accelerating delivery and reducing risk.

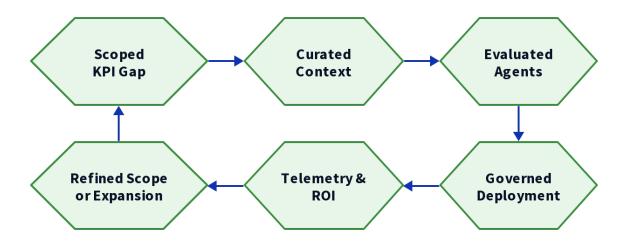


Figure 1.2: The measurable flywheel: each iteration compounds value and organizational learning

The cycle begins with a **scoped KPI gap** - a specific, measurable problem to solve. Teams curate **focused context** from GOLD knowledge sets rather than indexing everything. They build **evaluated agents** with automated quality gates ensuring citations, cost controls, and performance thresholds. These agents enter **governed deployment** with policy enforcement, immutable audit trails, and rollback procedures. Production generates **telemetry and ROI data** - real metrics on quality, cost, user satisfaction, and business impact. This evidence informs the next iteration:





either **refining scope** to address gaps revealed by production usage, or **expanding** to adjacent use cases using proven components. Each cycle compounds organizational learning, reusable assets, and stakeholder confidence.

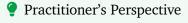
## 1.3 Reading guide

The chapters are structured to support both linear and selective reading:

- Chapter 2 Four-phase execution lifecycle (Scoping → Data → Modeling → Deployment) with practical guidance and phase-specific KPIs. Essential for implementation teams.
- Chapter 3 Six architectural viewpoints that capture decisions and enable governance. Includes the critical Governance, Risk & Compliance view (Section 3.2) and End-to-End Scenario validation (Section 3.6). Essential for architects and compliance teams.
- **Chapter 4** Comprehensive KPI catalog organized by lifecycle phase and architectural viewpoint, with telecom-specific metrics (MTTR, truck roll avoidance, FCR). Essential for defining baselines and success criteria.
- **Chapter 5** Practical use of the framework with real-life example of three-layer multi-agent architecture pattern, validated through production use cases (service recovery, network design, troubleshooting). Essential for solution architects and platform teams.
- Appendix A (Appendix) Classical ML workflow context for readers less familiar with data science practices, showing how traditional model development aligns with the agent lifecycle.

**Suggested reading paths:** Executives and governance stakeholders should focus on this introduction, Chapter 3, and Chapter 5. Implementation teams should follow the full sequence from Chapter 2 through Chapter 5.

Throughout the document, you will find practical advice and callouts.



Practitioner's Perspective callouts highlight practical, field-tested advice from real-world solution builders throughout this chapter. These sections distill lessons learned and best practices to help readers avoid common pitfalls and accelerate successful implementation.







# 2 Execution lifecycle

This chapter adapts the classical ML lifecycle (Appendix A) into a compact, action-oriented playbook for telcos. The goal is practical: achieve fast, controlled, and measurable AI outcomes while keeping the work reversible and auditable.

## 2.1 Adapted workflow at a glance

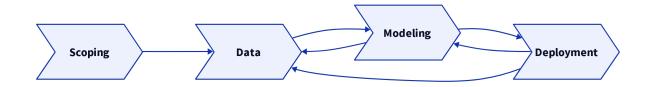


Figure 2.1: Lean agent system lifecycle showing iterative workflow phases with feedback loops

Many projects that combine AI and agents run into similar obstacles: unclear scope, overly ambitious data collection, fragile prototypes, and missing cost or quality governance. The workflow below is intentionally lean and reversible:

- 1. Scope tightly pick a single process and define unambiguous success metrics.
- 2. Curate a small, trusted knowledge core prioritize quality over raw volume.
- 3. Design and validate the "brain" (prompts, tools, and evaluation) before investing in large-scale tuning.
- 4. Deploy with guardrails monitor quality, latency, and cost, and ensure the ability to roll back within minutes.





Feedback arrows in the diagram indicate continuous improvement: production issues (slower-than-expected responses, missing facts, or rising inference cost, for example) should trigger targeted adjustments in earlier phases rather than broad rewrites.

## 2.2 Phase 1: Scoping (Project definition)

Objective: Identify a sharply defined, low-friction win that demonstrably proves value.

#### 2.2.1 Core activities

- Domain & use case selection pick a high-impact process that is low in hidden complexity. Deliverables should include precise inputs and outputs, a clear rule for when the agent should act instead of a human, and identified failure modes with fallback actions.
- Agent design start with a single agent for linear workflows. Introduce specialist or validator
  agents and conventional ML components only when differences in data type (for example,
  time-series KPIs versus tabular SLA metrics or images), complexity, or scale make them
  necessary.
- Stakeholder alignment map value to roles (RACI), set communication expectations, and curtail unmanaged "shadow" solutions early.
- Success criteria & KPIs choose conservative pilot targets (roughly 80% of current manual quality, for example). Capture user, business, and compliance metrics and define a closed feedback loop.

## 2.2.2 Architecture snapshot

#### Roles

- Foundation data and models (LLMs and, where appropriate, classical or specialist ML such as classification, clustering, forecasting, or anomaly detection), embeddings, and a vector database. This layer is responsible for security, storage, and lifecycle management.
- Orchestration executes agents and workflows; enforces validation, fact checks, and routing rules.
- Integration offers APIs and messaging to users and internal systems; implements access control and rate limiting.

## 2.2.3 Risks and mitigation

- Data privacy apply encryption and strict access controls in the Foundation; redact or mask sensitive fields at the API layer.
- Model drift run scheduled quality tests and version models and prompts consistently.
- Vendor lock-in keep an abstraction layer so providers can be swapped with minimal disruption.









**Figure 2.2:** Three-layer architecture showing Foundation (data/models), Orchestration (agents/workflow), and Integration (APIs/interfaces)

- Hallucinations introduce explicit validator or fact-check steps and clear human escalation paths.
- Governance maintain audit trails, security event logs, and records of model behavior.

Leading KPIs: user engagement, latency, data freshness

Lagging KPIs: efficiency gains, quality improvements, financial impact

## 2.2.4 Practical tips

- 1. Start with a well-understood, "boring" process.
- 2. Instrument from day one with traces, structured logs, metrics, and alerts.
- 3. Design for failure: explicit fallbacks and a human review path.
- 4. Communicate scope, timeline, and success metrics up front.

## 2.3 Phase 2: Data (Knowledge & context)

Objective: Provide only the reliable context required to answer real questions - no more.

The temptation at this stage is to index everything: every runbook, every ticket, every network diagram. In practice, more data often means more noise: outdated procedures, conflicting





documentation, and stale network configurations that confuse rather than clarify. A tightly curated knowledge base of 20 validated documents will consistently outperform a sprawling collection of 2,000 unvetted files. Quality beats volume every time, and your operational staff will thank you when the system delivers precise answers instead of plausible-sounding hallucinations.

NOTE: If the use case includes structured telemetry, time-series counters, ticket lifecycle events, or predictive risk scores, plan early how classical ML outputs (forecasts, anomaly flags, churn or failure probabilities) will be presented as discrete "facts" to the LLM rather than as raw feature dumps. Maintain a simple contract: ML service  $\rightarrow$  feature synthesis  $\rightarrow$  compact retrieval context.

#### 2.3.1 Core activities

- Source & curate identify authoritative systems and tier data into GOLD (human-checked), SILVER (historically reliable), and BRONZE (raw or synthetic). Begin with a small GOLD set (on the order of 100–500 examples).
- Retrieval (RAG) design set chunking rules, pick an initial embedding model shortlist, and define the promotion path (dev → prod → enterprise). Measure before you expand.
- API/context integration add circuit breakers, short-term caching (for example, 5–15 minutes), and rate limits; surface freshness timestamps to users.
- Workflow mapping verify actual process steps and enumerate decisions, approvals, and exceptions.
- Governance version datasets, capture lineage, classify access levels, and define retention policies.
- Monitoring measure relevance, latency, freshness, drift, and cost per retrieval. For numerical or time-series feeds, also track data lag, percent missing, and anomaly flag precision.

## 2.3.2 Design choices and rationale

- Dataset tiers reserve GOLD for human-validated material used in evaluation; add lower tiers only once the process is stable.
- Embedding model start with a solid off-the-shelf model and defer custom tuning until measured metrics indicate a persistent deficit.
- Vector store progression use lightweight stores (FAISS/Chroma) in development, migrate to managed stores (Qdrant/Pinecone/Weaviate) in production, and consider enterprise search only if governance or scale requires it.
- Retrieval contract cap tokens per query (aim for roughly 4–8K) and introduce re-ranking only if precision degrades.
- Resilience baseline combine gateway, circuit breaker, cache, and rate limiter as a reusable template.





 Data quality monitoring - automatic anomaly and drift checks should trigger alerts and quarantine. Treat drift detection for structured features separately from drift in text embeddings.

## 2.3.3 Risks and mitigation

- Stale or incorrect context enforce freshness SLAs, display timestamps, and skip stale content.
- Privacy leakage classify and redact PII early and enforce least-privilege access controls.
- Rising cost monitor tokens per answer and remove unused or low-value embeddings.
- Over-engineering begin with a single retriever and add hybrid search only if recall is demonstrably insufficient.

Leading KPIs: freshness percentage, retrieval precision, average tokens retrieved, drift score, and cost per query.

Lagging KPIs: answer correction rate, escalations caused by missing context, and user trust scores.

## 2.3.4 Practical tips

- 1. Ship a thin vertical slice: one authoritative source  $\rightarrow$  one embedding model  $\rightarrow$  one vector DB.
- 2. Version all parameters (dataset, embedding model, chunking) for auditability.
- 3. Measure before tuning: keep a small GOLD Q&A set for regular precision/recall checks.
- 4. Automate hygiene: nightly jobs to remove stale or orphaned embeddings.

## Practitioner's Perspective

By designing your data layer with this evolution path in mind - simple tools first, enterprise replacements later - you avoid the painful "rip and replace" scenarios that have derailed many proof-of-concepts. Keep interfaces abstract enough that swapping vector stores requires configuration changes, not code rewrites.

## 2.4 Phase 3: Modeling (Agent design, ML integration & evaluation)

Objective: Assemble prompts, tools, and agent logic, and validate quality, latency, and cost before committing to fine-tuning.

Many teams prematurely jump to fine-tuning or complex multi-agent architectures. The counterintuitive reality: you'll likely achieve 80-90% of your target quality through careful prompt engineering and smart tool design. Fine-tuning should be your last resort, not your first move-it locks you into a specific model version, requires ongoing maintenance, and costs significantly more than prompt iteration. Start simple, measure obsessively, and only add complexity when data proves it's necessary.





#### 2.4.1 Core activities

- Base model selection choose a model that matches the task (reasoning depth, latency, cost, and any required security guarantees). Begin with cost-effective options and upgrade only if measured gaps persist.
- Prompt template keep a stable structure (role, context, task, constraints, examples, output format) and version it.
- Evaluation harness combine automated tests with periodic human spot checks; include regression and A/B comparisons.
- Tool integration define strict JSON schemas, label tools by expected latency (fast <3s, slow >3s, human-gate, fallback), and specify error handling. Surface predictive ML endpoints (forecast\_next\_outage, classify\_ticket\_priority, detect\_anomaly, for example) as explicit, versioned tools with narrow schemas.
- Multi-agent split introduce multiple agents only when specialization or parallelism demonstrably improves outcomes.
- Safety & quality add hallucination detection, format validators, and compliance keyword scans.

## 2.4.2 Design choices and rationale

- Model portfolio maintain a primary model and a fallback (different vendor or modality) to improve resilience.
- Prompt governance keep prompts under version control, use a review process, and log prompt hashes per response.
- Tool call contract enforce strict schema validation, retries, and circuit breakers for unreliable APIs.
- State & memory separate short-term conversation state from longer-term storage (vector DB or database), and apply expiry rules and size limits. Treat ephemeral reasoning traces differently from durable ML indicators (a rolling risk score, for instance) to prevent memory bloat.
- Evolution path begin with a single agent, add tools, and only move to a specialist trio (Orchestrator / Domain / Validator) when metrics justify the added complexity.

## 2.4.3 Risks and mitigation

- Prompt drift require formal reviews and trigger alerts on unapproved changes.
- Latency spikes limit parallel tool calls and cache deterministic data where appropriate.
- Cost overrun monitor tokens per turn and automatically switch to a cheaper model for low-risk queries.
- Inconsistent output apply strict format checks and automated repair with retries when needed.
- Hallucinations require cited sources and use a validator agent for fact checking.





Model KPIs: quality score, citation coverage percentage, format compliance rate, P95 latency, cost per successful task, and hallucination rate. When predictive ML tools are present: model precision and recall (or ROC AUC), feature drift score, forecast MAPE/RMSE, and anomaly alert precision and recall.

## 2.4.4 Practical tips

- 1. Extract gains from prompt and tool design (including ML tool calibration) before considering fine-tuning fine-tuning typically requires on the order of 1K high-quality examples to be cost-effective.
- 2. Automate evaluation early to prevent hidden regressions.
- 3. Keep the tool list short each tool adds complexity and a potential source of drift; retire those with low value quickly.
- 4. Plan for rollback: preserve the ability to replay requests with the exact model, prompt, and toolset (including ML model versions).
- 5. Separate evaluation layers: LLM answer quality, predictive model performance, and retrieval relevance.

## Practitioner's Perspective

The industry buzz around multi-agent systems can be misleading. Specialized agents handling distinct domains can improve accuracy and reduce latency through parallelism, that's true. However, they also multiply your debugging surface, complicate your evaluation strategy, and increase operational overhead. In telecom use cases, we've found that most pilots succeed with a single well-designed agent. Only introduce multiple agents when you have clear evidence - not assumptions - that specialization solves a measurable problem.

## 2.5 Phase 4: Deployment (Production, ML Ops and monitoring)

Objective: Operate agents safely in production with full visibility and the ability to roll back quickly.

Deployment is not a finish line, it's the beginning of continuous learning. Use real user interactions to validate assumptions, surface edge cases, and guide iterative improvement. This requires a mindset shift: instead of trying to perfect the system before launch, you ship something good enough with comprehensive instrumentation, then evolve based on evidence. The key enabler is reversibility. Every change must be undoable within minutes, not hours or days.

#### 2.5.1 Core activities

• Topology & placement - determine where sensitive and general data should reside (for example, keep PII on-prem, run inference in a private cloud, and cache at the edge only where safe).





- Interface exposure present a single API gateway that enforces authentication, rate limits, and schemas.
- Identity & access use single sign-on, role or attribute-based access, regular key rotation, and per-request policy checks.
- Observability capture traces, structured logs, and both technical and business metrics; tag each call with model and prompt versions and, where relevant, ML model version and feature bundle ID.
- Rollout & incidents adopt phased rollouts, canary testing, and clear rollback playbooks.
- Cost & compliance apply cost tagging for chargeback, dashboards, audit logging, and retention/deletion jobs. Attribute ML inference cost separately from LLM token spend for optimization.

## 2.5.2 Design choices and rationale

- Deployment topology avoid unnecessary movement of sensitive data and, where possible, bring compute to the data. Co-locate high-volume predictive services with their data sources to cut latency and egress.
- Gateway contract a single front door that rejects unsafe or malformed requests.
- SLOs for example, aim for P95 latency under 3s, 99.9% uptime, and agent error rates below 0.1%; tie alerts to business impact thresholds.
- Auditability capture inputs, retrieved context IDs, model and prompt versions, tool calls, outputs, user IDs, and ML model versions with input feature hashes.
- Rollout strategy internal beta → pilot group → percentage rollout → full deployment only after guardrail checks pass.

#### 2.5.3 Risks and mitigation

- Cost growth set budgets and alerts and automatically switch to a cheaper model for low-risk traffic spikes.
- Security gaps perform secret scanning, enforce least-privilege, and adopt deny-by-default policies.
- Silent quality drops gate promotions with continuous evaluation.
- Operational toil automate runbooks and add self-healing procedures where practical.

Operational KPIs: uptime, P95 latency, error rate, rollback count, cost per 1K requests, MTTR. When ML tools are present: ML inference latency, feature pipeline freshness, model drift alerts count.

Business KPIs: adoption %, automation rate, escalation frequency, time-to-value.





#### 2.5.4 Practical tips

- 1. Gate releases with both SLO and quality checks (include health gates for predictive ML models).
- 2. Keep changes reversible: feature flags, shadow traffic, and quick rollback (under 5 minutes or so) for both LLM prompt versions and ML model versions.
- 3. Unify telemetry: use correlation IDs across components and propagate model, prompt, and ML version IDs.
- 4. Track cost by use case or tenant to enable optimization (separately attribute LLM token spend and ML inference cost).
- 5. Schedule regular ML evaluation jobs alongside LLM quality regression checks.

## Practitioner's Perspective

The audit trail deserves special attention in telecom contexts. Regulatory requirements, customer disputes, and internal quality reviews all demand detailed records of how the system reached specific decisions. However, comprehensive logging isn't just about compliance - it's your debugging superpower. When users report unexpected answers, you need to replay the exact retrieval results, model version, and reasoning chain that produced that response. Tag everything, version everything, and retain logs longer than you think necessary.

## 2.6 Putting it all together

The loop (Scope  $\rightarrow$  Data  $\rightarrow$  Modeling  $\rightarrow$  Deployment) becomes a productive flywheel when each phase hands off concrete artifacts to the next: scoped KPIs focus data curation; curated GOLD data supports reliable evaluation; evaluation results steer prompt, tool, and model updates; and production telemetry validates or challenges earlier assumptions, guiding targeted fixes rather than broad rewrites.

#### A compact playbook:

- Start narrow: capture a one-sentence use case and 3–5 KPIs.
- Build a small GOLD corpus and measure retrieval performance before expanding.
- Treat the prompt and tool stack as the primary optimization surface before investing in tuning.
- Instrument everything from the prototype stage (trace IDs, versioning, cost metrics).
- Make reversibility non-negotiable: feature flags, rollbacks, and version logs.

Practical success typically hinges less on elaborate agent graphs and more on disciplined iteration with observable, governed change. When conventional ML models are part of the system, treat them as first-class, versioned tools - audited, evaluated, and charged separately. Keep the system lean, measurable, and reversible; scale only what metrics justify.





Next: Chapter 3 converts this workflow into a lightweight architectural viewpoints framework. Each viewpoint (lifecycle, governance, runtime, capability, infrastructure, scenario) records just enough decisions, artifacts, and KPIs to guide evolution without bogging teams down - and provides an auditable rationale for changes made across Scope  $\rightarrow$  Data  $\rightarrow$  Modeling  $\rightarrow$  Deployment.





# 3 Architectural viewpoints

This chapter adapts the time-tested architectural viewpoint approach - originally developed for complex distributed systems - to the specific challenges of AI agents in telecom operations. Rather than comprehensive upfront documentation, we apply the viewpoint discipline selectively: capturing decisions where they matter most, versioning them alongside code, and retiring outdated entries as the system matures.

Viewpoints structure decisions that otherwise scatter across slides and chat threads. They are templates or specifications that are used to describe particular aspects of a system's architecture and provide focused lenses for different stakeholders to understand, evaluate, and govern the system effectively.

Views are concrete instantiations of those viewpoints for a specific system, populated with actual data, configurations, and decisions. Think of them as compact, reviewable lenses you update as the system evolves - each records why a choice was made, what evidence supported it, and what to monitor next. You rarely need every view at full depth for a telecom use case: start with the few that address real risk or unblock delivery, and expand only when metrics show ambiguity, drift, or scaling pressure.

In practical terms, think of views as living documents that grow with your system. A newly scoped pilot might have a two-page Data & Model Lifecycle view and a sketch for Runtime topology. Six months into production, those same views might span ten pages each, enriched with operational learnings, incident post-mortems, and optimization decisions. The goal isn't comprehensiveness from day one - it's capturing just enough structure to enable informed decision-making without premature over-design.

## Why they matter:

- Focused lenses each viewpoint isolates one concern (data lineage, risk controls, execution graph, capability roadmap, infrastructure, or the concrete user journey), preventing a single diagram from swallowing every decision.
- Role targeting network operations focus on runtime failure domains and SLOs; compliance inspects governance logs; ML engineers set lifecycle drift thresholds; product owners validate scenarios; platform teams evolve infrastructure knobs.
- Cross-link traceability the Scenario view serves as the living thread: every tool call, retrieval rule, latency budget, and policy decision that appears there should map back to an entry in one of the other five catalogs.







- Early inconsistency surfacing mismatches (a latency budget in Runtime that contradicts a P95 target in Physical, or a retention policy that clashes with a feature freshness window, for instance) become tracked deltas rather than hidden technical debt.
- Lean governance decisions are versioned where they belong; diffs on a view record architectural change without digging through source or chat logs.

The end-to-end Scenario view is deliberately last: it pressure-tests whether each observed behavior in a real request can be traced to a concrete, versioned entry in the other viewpoints.

This chapter maps the iterative workflow (Scope  $\rightarrow$  Data  $\rightarrow$  Modeling  $\rightarrow$  Deployment) into six compact architectural viewpoints.

Goal: capture only those decisions that accelerate safe progress, enable auditability, and avoid over-design. The items inside a view typically emerge from phase work, but the viewpoints themselves act as orthogonal catalogs of decisions and artifacts you can version, review, and reuse.

## 3.1 Data & Model Lifecycle viewpoint

Purpose: Trustworthy, evolving knowledge and model assets that feed Modeling and Deployment.

A view that instantiates this viewpoint for a particular system answers the fundamental question: "Where does our AI get its knowledge, and how do we know it's reliable?" In telecom environments with multiple documentation systems, legacy databases, and continuously changing network configurations, maintaining data quality is an ongoing operational challenge, not a one-time setup task.

#### Capture:

- Dataset tiers (GOLD/SILVER/BRONZE), lineage, retention, and PII handling.
- Feature and embedding generation configuration (chunking, embedding model id, refresh cadence).
- Model registry entries (model id, data hash, hyperparameters, evaluation set versions).
- RAG contract max tokens retrieved, top-k, re-ranking policy, and provenance rules.
- Quality gates data constraints, drift thresholds, and golden-set metrics.

## Example: GOLD Dataset Entry

Dataset: network-topology-docs-v3

Tier: GOLD

Source: Confluence wiki + validated network diagrams

Lineage: Manual export → JSON transform → validation → embedding

Refresh: Weekly (Sundays 02:00 UTC)





```
Embedding: text-embedding-ada-002, chunk_size=512, overlap=50 Retention: 90 days (rolling), PII: None (network data only)
```

Quality gate: Min 95% retrieval precision on golden Q&A set (n=50)

#### KPIs:

- Freshness percentage (% of data within SLA refresh window)
- Retrieval precision (relevant docs in top-k results)
- Drift score (embedding distribution shift week-over-week)
- Model evaluation metrics (MAPE, precision, hallucination rate for cited answers)

Phase inputs/outputs: Data phase activities produce the GOLD seeds and retrieval configuration captured here; Modeling consumes registry entries and returns evaluation findings; Deployment telemetry updates drift and freshness indicators logged in this view.

## 3.2 Governance, Risk, and Compliance viewpoint

Purpose: Controlled, explainable evolution that reduces regulatory and reputational risk across phases.

Governance is where AI ambitions meet corporate reality. A view that materializes this viewpoint documents not just what policies exist, but how they're enforced in code, monitored in production, and audited after incidents. For telecom organizations subject to GDPR, industry regulations, and internal compliance frameworks, this viewpoint often becomes the most scrutinized artifact during audit reviews and executive briefings.

#### Capture:

- Policies (safety, privacy, IP) and enforcement points (pre/post filters, tool gating).
- Human oversight flows (escalation triggers and review SLAs).
- Audit log schema (prompt id, model version, tool calls, retrieved doc ids, decision codes).
- Incident playbooks (hallucination spike, data leak, model regression).

## Example: PII Redaction Policy

Policy: PII-redaction-001

Enforcement: Pre-processing filter on all user inputs Rule: Mask phone numbers, email addresses, credit cards

Audit: Log original hash, redacted version, timestamp, user\_id Escalation: Security team notified if >10 PII detections/hour Review SLA: Weekly audit of flagged content (compliance team)

Incident playbook: PII-LEAK-2025 (lock system, notify DPO, forensics)





- Policy violation rate (violations per 1K requests)
- False-block rate (legitimate requests blocked by policies)
- Review turnaround time (hours from escalation to resolution)
- Incident MTTR (mean time to resolution for governance incidents)

Phase inputs/outputs: Scoping sets initial risk boundaries referenced here; Data stewardship records lineage used in audits; Modeling adds new evaluation suites that extend policy tests; Deployment supplies enforcement and incident metrics.

## 3.3 Runtime/Process (Execution Topologies) viewpoint

Purpose: Reliable, debuggable execution graph for agents and ML tools.

When things go wrong in production, an instance of this viewpoint is your diagnostic map for the system you build. It shows how a user request flows through your system, where decisions happen, what can fail, and how failures are handled. The difference between a 5-minute incident resolution and a 5-hour debugging session often comes down to how well you've documented your runtime topology.

#### Capture:

- Chosen topology (single agent, orchestrator-worker, hierarchical planner) and rationale.
- Explicit node list: retrieval, routers, predictive ML tools, validators, and human gate.
- Failure handling (timeouts, retries, circuit breakers, fallback model escalation).
- Cost and latency budgets per node, plus caching and parallelism rules.

## Example: Orchestrator-Worker Topology

Topology: Orchestrator-worker (3 domain agents)

#### Nodes:

- Orchestrator: Intent classification (GPT-4o-mini, <500ms)
- Network-ops-agent: Troubleshooting domain (GPT-4o, <3s)
- RAG-retriever: Knowledge Graph query (<200ms)
- Validator: Policy compliance check (<100ms)

#### Failure handling:

- Retrieval timeout: 5s → fallback to cached results
- LLM timeout: 30s → escalate to human
- Circuit breaker: Open after 3 consecutive LLM failures

Cost budget: \$0.05 per request (P95), Latency: 3s end-to-end (P95)

- P95 latency (overall and per critical node)
- Tool failure rate (failed invocations per 1K requests)





- Escalation frequency (human intervention rate)
- Cost per successful task (dollars per completed request)

Phase inputs/outputs: Scoping limits scope and complexity, which informs the node set; Data supplies retrieval endpoints; Modeling produces the initial graph and tool contracts; Deployment provides live latency and failure statistics for pruning.

## 3.4 Capability/Development viewpoint

Purpose: Platform capability maturity roadmap so teams invest in the right enablers at the right time.

This viewpoint prevents the common trap of building infrastructure before you understand requirements. Instead of guessing what tools and platforms you'll eventually need, you track what capabilities your current use cases actually require and plan incremental investments. It's the antidote to both under-investment (brittle point solutions) and over-investment (unused enterprise platforms).

#### Capture:

- Current versus planned capabilities (orchestration, evaluation harness, tool registry, model gateway, feature store, router confidence calibration).
- Build versus buy decisions and modular interfaces (tool schema, retrieval API, model contract including confidence/abstain signals).
- Versioning and promotion workflow (prompt review cycle, tool schema change policy).

## Example: Evaluation Harness Build Decision

Capability: Evaluation harness

Status: Build (in-house)

Rationale: Off-the-shelf tools lack telecom domain metrics

Components:

- Golden Q&A set (n=200, telecom-specific scenarios)
- Automated regression suite (runs on every prompt change)
- Human eval UI for edge cases (weekly review sessions)

Maturity: Level 3 (automated, versioned, production-integrated)

Roadmap: Add A/B testing framework (Q2 2026)

Interface: REST API (/evaluate endpoint, JSON schema v2.1)

- Lead time for prompt or tool change (hours from request to production)
- Regression escape rate (bugs reaching production per release)
- Evaluation coverage percentage (% of use cases with golden sets)





• Reusable component adoption (% of new use cases reusing existing tools)

Phase inputs/outputs: Early Scoping highlights capability gaps; Data and Modeling iterations add concrete components (retriever variants, evaluation harness); Deployment feedback reprioritizes the roadmap captured here.

## 3.5 Physical/Infrastructure viewpoint

Purpose: Performance, resilience, and cost controls that support Deployment and feed back operational telemetry.

Infrastructure decisions cascade through the entire system lifecycle. This viewpoint allows documenting the concrete deployment topology - where models run, how data flows, what breaks under load, and who pays for what. In telecom contexts with strict uptime requirements and cost constraints, these choices directly impact both operational success and business viability.

#### Capture:

- Serving layers (model gateway, batching, cache tiers) and scaling strategy (HPA metrics, GPU pool segmentation).
- Storage tiers (vector DB versus object store), backup, and immutability.
- Network and security (egress policy, private endpoints, secret management, rate limits at gateway).
- Observability stack (traces tagged with model/prompt/tool versions, and redaction rules).
- Cost governance (per-tenant budgets, token versus ML inference breakdown).

## Example: Cloud Infrastructure Configuration

Serving: Azure OpenAI + self-hosted embedding service Scaling: HPA on request queue depth (target: 100 pending) Storage:

- Vector DB: Qdrant (3-node cluster, 1TB SSD)
- Object store: Azure Blob (GOLD datasets, immutable after validation)
- Backup: Daily snapshots, 30-day retention

#### Security:

- Private endpoints only (no public internet)
- Secrets: Azure Key Vault with rotation every 90 days
- Rate limits: 100 req/min per tenant (burst: 150)

Observability: OpenTelemetry → Azure Monitor, traces include model\_version tag Cost: \$0.02/1K requests (target), \$3K/month budget cap

#### KPIs:

• Uptime (SLA: 99.9% monthly)





- P95 latency (target: <3s end-to-end)
- Error rate (4xx + 5xx errors per 1K requests)
- Cache hit percentage (target: >60% for retrieval queries)
- Cost per 1K requests (dollars, tracked by use case)
- Rollback count (production rollbacks per month)

Phase inputs/outputs: Scoping defines SLO baselines; Data ingestion choices drive storage and freshness tactics; Modeling determines serving needs (model size, tool mix); Deployment operational metrics populate such views over time.

## 3.6 End-to-End Scenario viewpoint

Purpose: A validated trace of a representative request that shows guardrails, data usage, and decision points - evidence that the prior views connect coherently.

The Scenario view is where abstract architecture meets concrete reality. It captures real end-to-end request flows - successful and failed-showing how governance policies trigger, which data sources contribute, how models make decisions, and where costs accumulate. This view serves both as validation (do all the viewpoints actually work together?) and as debugging reference (what happened in production?).

## Capture:

- Sequence diagram with Gateway → Orchestrator → Policy → Tools/Retrieval → Models →
  Post-checks → Monitoring.
- Annotated decisions (router scores, thresholds, abstain logic, tool schemas used).
- Trace ID example with linked logs and evaluation outcomes.

## **i** Example: Network Troubleshooting Trace

Scenario: Network troubleshooting request

Trace ID: trace-2025-10-29-abc123

#### Flow:

- 1. Gateway: Auth (JWT valid) → rate limit check (pass)
- 2. Orchestrator: Intent="network-ops" (confidence: 0.94)
- 3. Policy: PII check (none detected) → compliance pass
- 4. Retrieval: Query Knowledge Graph (3 docs, 2.1K tokens)
- 5. Network-ops-agent: GPT-4o reasoning (15 steps, 8.2s)
- 6. Tools: SDN API call (link status query, 340ms)
- 7. Post-check: Validator confirms no hallucinations (0.97 groundedness)
- 8. Monitoring: Log trace + cost (\$0.04) + latency (9.1s total)

Outcome: Success, user satisfied (feedback: thumbs-up)





- Scenario task success rate (% of requests completed successfully)
- Groundedness score (avg citation quality, 0-1 scale)
- Average decision path length (number of tool invocations per request)
- Human intervention rate (% of requests requiring escalation)

Phase inputs/outputs: Scenario steps are assembled once initial Scoping, Data, and Modeling artifacts exist; Deployment traces enrich the scenario with real latency, tool path variance, and grounding metrics for regression guarding.

## 3.7 Applying viewpoints incrementally

The six viewpoints are best adopted incrementally, aligned with the Scope  $\rightarrow$  Data  $\rightarrow$  Modeling  $\rightarrow$  Deployment workflow (Chapter 2). Start with a minimal viable set, adding depth only when real risks or bottlenecks emerge.

## 3.7.1 Adoption progression

## Initial focus (Scoping & Data phases):

Begin with three foundational viewpoints that form the core feedback loop:

- Data & Model Lifecycle: Capture GOLD seeds, retrieval contracts, and quality gates. This records which data and models drive decisions.
- Runtime/Process: Document execution topology, initial node list, and latency budgets. This shows how artifacts flow through the system.
- **End-to-End Scenario:** Define the target request flow structure. While full trace evidence only arrives during Deployment, the *template* (expected sequence, decision points, validation criteria) guides Data and Modeling choices.

Why these three: Scoping sets boundaries; Data phase produces GOLD seeds and retrieval contracts that populate the Lifecycle view and anchor Runtime decisions. The Scenario template provides early validation criteria without premature complexity.

## Pre-deployment hardening (late Modeling, pre-Deployment):

Add governance and infrastructure viewpoints as the system approaches production:

- Governance, Risk & Compliance: Record policy boundaries, enforcement points, audit schema, and incident playbooks. As Modeling finalizes tool schemas and evaluation suites, this view captures where and how policies are enforced.
- Physical/Infrastructure: Document latency budgets, scaling triggers, observability configuration, and cost governance rules. Deployment planning exposes infrastructure needs (serving layers, storage tiers, monitoring) before live traffic.





Why now: These viewpoints provide explicit guardrails and operational controls. Adding them too early creates unused overhead; adding them too late risks compliance gaps or performance surprises in production.

#### Operational maturity (Deployment & beyond):

Expand capability tracking based on production experience:

• Capability/Development: Track build-vs-buy decisions, feature maturity levels, and platform roadmap driven by operational feedback. Early iterations reveal gaps (e.g., router confidence calibration, feature store requirements, evaluation harness improvements).

Why last: Premature capability planning creates infrastructure before understanding requirements. Operational reality - which tools see heavy use, which integrations bottleneck, which governance policies trigger frequently - should drive capability investment.

## **Governance throughout:**

Treat each view as a versioned set of documents with change records capturing: - What changed (e.g., added fraud detection policy to Governance view) - Why it changed (e.g., hallucination spike in production telemetry) - Metrics before/after (e.g., false-positive rate, review latency)

This preserves auditability and surfaces architectural drift. Reviewers can diff view versions rather than reconstructing decisions from source code or chat logs.

## 3.7.2 Viewpoint adoption by phase

The table below shows which viewpoints to enable at each lifecycle phase and what concrete artifacts become part of each view. Read it as your incremental adoption roadmap: start with three viewpoints (Data & Model Lifecycle, Runtime, Scenario) during Scoping & Data, progressively layer in Governance and Physical constraints as you approach production, and finally add Capability tracking once operational patterns emerge. The "↑ (above)" notation indicates cumulative adoption - each phase builds on the viewpoints from previous phases.

**Table 3.1:** Viewpoint adoption progression aligned with lifecycle phases

Phase	Active Viewpoints	Key Outputs
Scoping & Data	Data & Model Lifecycle, Runtime topology, Scenario template	GOLD seeds, retrieval contract, execution graph template, expected trace structure
Modeling	↑ (above) + Governance (start drafting)	Evaluation suites, tool schemas, initial policy boundaries
Pre- Deployment	↑ (above) + Governance, Physical	Audit schema, SLO targets, infrastructure config





Phase	Active Viewpoints	Key Outputs
Deployment & ops	All six viewpoints	Live telemetry, incident playbooks, real drift metrics, actual trace evidence

When to expand: Add view depth when data lineage complexity emerges, tool integration patterns stabilize, latency/cost constraints tighten, or compliance mandates surface. The roadmap evolves continuously based on operational feedback.

## 3.7.3 View-to-phase integration

The previous table showed *when* to activate viewpoints; this table reveals *how* they expand with depth across phases. Read each row to understand a single viewpoint's lifecycle: which phases contribute to the view (columns 2-5) and which viewpoints use them (rightmost "Feeds To" column). For example, the Data & Model Lifecycle row shows it receives GOLD seeds during the Data phase, evaluation findings during Modeling, and drift metrics during Deployment - then feeds lineage information to Governance, endpoint configurations to Runtime, and freshness SLAs to Physical infrastructure planning. The "-" symbol indicates phases where the viewpoint is not likely updated with phase-specific artifacts (though the viewpoint may still be maintained or referenced).

Table 3.2: Viewpoint-to-phase integration matrix showing artifact flows

Viewpoint	Scoping	Data	Modeling	Deployment	Feeds To
Data & Model Lifecycle	-	GOLD seeds, retrieval config	Evaluation findings, model registry	Drift metrics, freshness telemetry	Governance (lineage), Runtime (endpoints), Physical (freshness SLAs)
Governance, Risk & Compliance	Risk boundaries	-	Policy tests	Enforcement metrics, incidents	Runtime (policy hooks), Scenario (compliance traces)





Viewpoint	Scoping	Data	Modeling	Deployment	Feeds To
Runtime/Process	Scope limits	Re- trieval end- points	Tool contracts, execution graph	Latency stats, failure rates	Physical (SLOs), Scenario (sequence structure)
Capability/Development	Gap inventory	Compo- nents	Maturity tracking	Operational feedback	Platform investment roadmap
Physical/In- frastructure	SLO baselines	Volume, fresh- ness needs	Serving needs (model size, tool mix)	Actual uptime, cost per request	Runtime (failure modes), Scenario (latency traces)
End-to-End Scenario	Expected flow template	-	-	Actual trace evidence	Regression tests, review checklists, readiness criteria

Result: A lightweight, evolving architecture narrative that keeps decisions transparent while preserving the rapid iteration loop established in Chapter 2. Views grow with the system, capturing just enough structure to enable informed decision-making without premature over-design.





# 4 Key metrics & KPIs

This chapter consolidates metrics from the four-phase lifecycle (Chapter 2) and six architectural viewpoints (Chapter 3) into a practical measurement framework. The goal: instrument agents and ML tools with metrics that enable fast, evidence-based decisions while maintaining governance and cost control.

Treat this as a curated catalog of high-impact metrics proven effective in real-world telecom deployments - not an exhaustive list, but a starting point to adapt for your specific use case. When defining your own metrics, be specific: provide measurable definitions, baseline measurements, target thresholds, and clear success criteria that reflect both business impact and technical health. The values shown are illustrative examples; set targets relevant to your operational context.

## 4.1 Metrics by lifecycle phase

This table organizes key metrics according to the Scope  $\rightarrow$  Data  $\rightarrow$  Modeling  $\rightarrow$  Deployment workflow, aligned with metrics defined in Chapter 2.

**Table 4.1:** Lifecycle phase metrics summary

Phase	Leading KPIs (predict progress)	Lagging KPIs (validate value)	Example Targets
Scop- ing	User engagement rate, P95 latency, Data freshness %	Efficiency gains, Quality improvements, Financial impact	40% adoption in 6 weeks; <3s response time; Positive ROI in 3 months
Data	Freshness %; Retrieval precision; Avg tokens retrieved; Drift score; Cost per query	Answer correction rate; Escalations (missing context); User trust scores	GOLD ≥95% fresh; ≥80% precision at top-5; <15% correction rate





Phase	Leading KPIs (predict progress)	Lagging KPIs (validate value)	Example Targets
Mod- eling	Quality score; Citation coverage %; Format compliance; P95 latency; Cost per task; Hallucination rate	ML model precision/recall/ROC AUC; Feature drift score; Forecast MAPE/RMSE; Anomaly alert precision	≥85% quality (pilot); ≥92% quality (production); <2% hallucination rate; ROC AUC ≥0.90
De- ploy- ment	Uptime; P95 latency; Error rate; Cost per 1K requests; MTTR; ML ops: ML inference latency, Feature pipeline freshness, Model drift alerts	Adoption %; Automation rate; Escalation frequency; Time-to-value	≥99.9% uptime; <\$2 per 1K requests; ≥60% adoption in 3 months; ML: <500ms (fast tools), <15min freshness

**How to measure:** Each metric requires defining data sources (API logs, user surveys, ticketing systems), computation methods (formulas, aggregation windows), and ownership (who monitors and acts on thresholds).

## 4.2 Metrics by architectural viewpoint

This table maps metrics to the six viewpoints from Chapter 3, showing which KPIs validate each architectural concern.

**Table 4.2:** Viewpoint metrics summary

Viewpoint	Key Metrics	Example Targets	Data Sources
Data & Model Lifecycle	Freshness %; Retrieval precision; Drift score; Model evaluation (MAPE, precision/recall)	GOLD ≥95% fresh; ≥80% precision; Drift alert at >0.15	Vector DB metadata; Evaluation harness; Drift detection service
Governance, Risk & Compliance	Policy violation rate; False-block rate; Review turnaround; Incident MTTR	<1% violations; <5% false blocks; <4h high-priority reviews	Policy enforcement logs; User feedback; Incident management





Viewpoint	Key Metrics	Example Targets	Data Sources
Runtime / Process	P95 latency (per node); Tool failure rate; Escalation frequency; Cost per task	Retrieval <500ms; LLM <2s; Total <5s; <2% tool failures	Distributed tracing (OpenTelemetry); Tool orchestration logs
Capability / Development	Lead time (prompt/tool change); Regression escape rate; Evaluation coverage %; Reusable component adoption	<5 days for prompts; <5% regressions; ≥80% test coverage	Version control; CI/CD pipeline; Test suite metadata
Physical / Infrastructure	Uptime; P95 latency; Error rate; Cache hit %; Cost per 1K requests; Rollback count	≥99.9% uptime; <0.5% errors; ≥40% embedding cache; ≤1 rollback/quarter	Health monitoring; APM tools; Cache layer metrics
End-to-End Scenario	Scenario task success rate; Groundedness score; Avg decision path length; Human intervention rate	≥95% success; ≥90% groundedness; 2-5 steps (simple); <15% intervention	Scenario test suite; Citation validator; Trace logs

## 4.3 Telecom-specific KPIs

These metrics address operational concerns unique to telecom environments, grounded in real Network Operations Center (NOC) and customer service workflows.

**Note:** Actual baseline performance and achievable targets vary significantly by operator, use case complexity, existing automation maturity, and organizational context. Measure your specific baseline before setting targets.

**Table 4.3:** Network operations KPIs

Metric	Definition	Example Baseline → Target	Data Source
MTTR reduction	Mean Time to Resolve network incidents	4.2h → 2.5h (40% reduction)	NOC ticketing system





Metric	Definition	Example Baseline → Target	Data Source
Truck roll avoidance	% field visits prevented by remote diagnosis	0% → 25% reduction	Field service management
First-time- right dispatch	% dispatches with correct parts/info first visit	65% → 85%	Technician completion reports
Network KPI correlation	% agent-identified anomalies matching real incidents	N/A → ≥75% accuracy	Anomaly detection logs + incident reports

**Cost impact example:** 25% truck roll reduction  $\times$  1,000 monthly dispatches  $\times$  \$200/truck roll = \$50,000/month savings.

Table 4.4: Customer operations KPIs

Metric	Definition	Example Baseline → Target	Data Source
First-contact resolution (FCR)	% issues resolved on first contact	55% → 70% (+15pp)	Call center analytics
Average handling time (AHT)	Mean time per customer interaction	8.5 min → 5.5 min (35% reduction)	Workforce management system
Service order error rate	% orders requiring manual correction	8-12% → <3%	Order management system
SLA compliance	% tasks completed within SLA	80-85% → ≥95%	Service assurance platform

**Business impact example:** 35% AHT reduction  $\times$  10,000 monthly calls  $\times$  3 minutes saved  $\times$  \$1/min labor = \$30,000/month savings.

## 4.4 Metric interdependencies and trade-offs

Understanding how metrics influence each other enables smarter optimization decisions.

## 4.4.1 Leading → Lagging relationships

These technical metrics predict business outcomes:





Leading MetricPredicts  $\rightarrow$  Lagging MetricResponse Time (P95) $\rightarrow$  User AdoptionRetrieval Precision $\rightarrow$  Answer Correction RateCitation Coverage $\rightarrow$  User Trust ScoreFormat Compliance $\rightarrow$  Downstream Error RateFreshness % $\rightarrow$  Escalation FrequencyTool Failure Rate $\rightarrow$  Automation Rate

**Table 4.5:** Leading-to-lagging metric relationships

**Application:** When leading metrics degrade, proactively fix root causes before lagging business metrics suffer.

#### 4.4.2 Common trade-offs

Optimizing one metric may negatively impact another:

- Accuracy vs. Latency: Higher accuracy often requires more retrieval, longer prompts, or validator steps (adds 1–3s latency).
  - *Mitigation:* Use fast models for simple queries; escalate to powerful models only when confidence is low.
- **Citation Coverage vs. Response Length:** Requiring citations increases output verbosity and token cost.
  - Mitigation: Compact citation format (footnote IDs instead of full URLs); cache frequently cited sources.
- Automation Rate vs. Quality: Aggressive automation may increase error rate if confidence thresholds are too low.
  - Mitigation: Calibrate confidence thresholds using ROC curves; set escalation trigger at optimal precision/recall balance.
- Cost vs. Freshness: Real-time data pipelines cost more than batch updates.
  - *Mitigation*: Tier data by sensitivity: real-time for alarms/outages, hourly for config, daily for documentation.
- Reusability vs. Specialization: Generic components enable reuse but may underperform vs. domain-tuned solutions.
  - *Mitigation:* Start generic; specialize only when KPI gaps persist after prompt and tool optimization.

## Practitioner's Perspective

Metrics are only valuable if they drive action. Establish clear ownership for each metric and define explicit thresholds that trigger investigation or rollback. A metric without an owner





and a response playbook is just noise. In telecom operations, where incidents can impact thousands of customers, the difference between a monitored metric and an actionable alert is often measured in minutes of downtime prevented.

## 4.5 Application guidance

For each metric, document:

- Metric name: Unique identifier (e.g., agent.retrieval.precision).
- **Definition:** Plain-English explanation and formula.
- Owner: Individual or team responsible for monitoring and action.
- Data source: Logs, APIs, or systems providing raw data.
- Target range: Acceptable thresholds and alert triggers.
- Review cadence: Daily monitoring, weekly review, or monthly deep dive.

## Practitioner's Perspective

Start by tracking fewer metrics well rather than many metrics poorly. Three accurately measured and actively monitored KPIs will drive better decisions than fifteen metrics that generate reports nobody reads. As your system matures and patterns emerge, expand your measurement coverage deliberately, adding new metrics only when they address specific questions or risks that have surfaced in operation.





# 5 Framework in action

The framework's value emerges through practical application: multi-agent architectures that scale across use cases while maintaining governance and measurability. This chapter demonstrates how the four-phase lifecycle (Chapter 2) and six architectural viewpoints (Chapter 3) guide concrete implementation decisions in telecom systems, validating the approach through real deployments and establishing reusable patterns for rapid replication.

### Practitioner's Perspective: Reusability

Real-world telecom deployments reveal a consistent architecture pattern: layered agent systems that translate the framework's abstractions into operational reality. Teams that master this architecture on their first use case can achieve 40-60% infrastructure reuse on the second, with potential acceleration to 70-80% reuse by the fifth deployment. That said, reuse percentages vary significantly based on domain similarity and organizational maturity.

### 5.1 Multi-agent reference architecture

Telecom agent systems converge on a three-layer architecture that embodies the framework's separation of concerns while enabling domain specialization and component reuse.

### **5.1.1** Layer 1: Orchestrator (intent classification)

The orchestrator serves as the system's entry point, routing user requests to appropriate domain specialists based on intent detection. This layer implements governance controls defined in views based on the Governance, Risk & Compliance viewpoint (Section 3.2) through authorization checks, compliance validation, and audit trail generation. It maps directly to the framework's **Scoping phase** (Section 2.2) - understanding request context, validating boundaries, and determining the appropriate execution path.

#### Key responsibilities:

- Intent classification: determine which domain agent handles the request
- Authorization enforcement: validate JWT claims against allowed operations
- Request routing: dispatch to specialized domain agents
- Audit logging: capture request metadata for compliance tracking







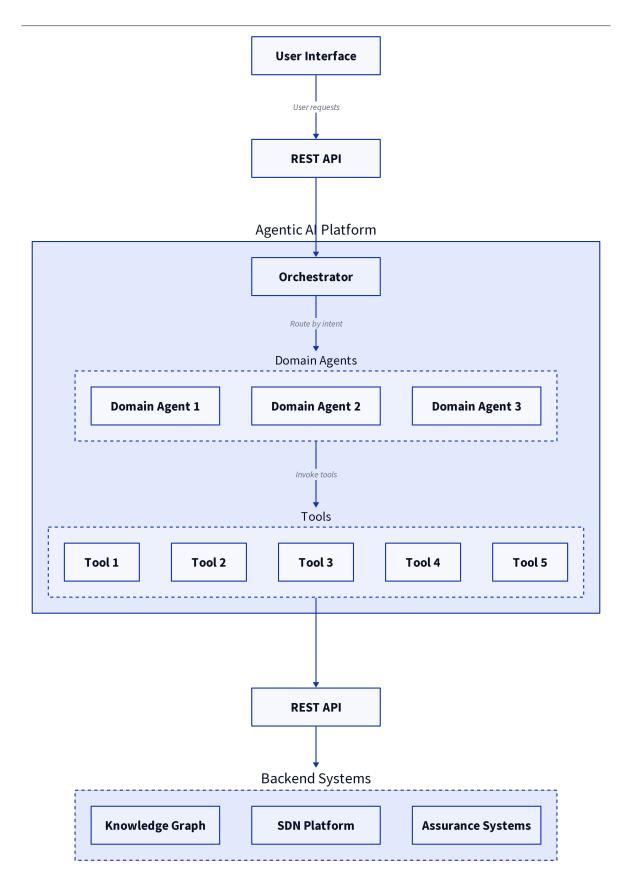


Figure 5.1: Three-Layer Agent Architecture for Telecom





**Solution examples:** LangGraph (multi-agent orchestration); Kong Gateway, Azure API Gateway (authentication, rate limiting, audit logging).

### 5.1.2 Layer 2: Domain agents (planning and execution)

Domain agents specialize by operational area: network operations, service assurance, provisioning, customer support. Each agent implements the full **four-phase lifecycle** (Chapter 2) internally, maintaining isolation between domains while sharing common infrastructure.

Within each domain agent:

- Scoping: Validate request against service rules and domain constraints
- Data: Gather context from Knowledge Graph, telemetry systems, and domain APIs
- Modeling: Apply ReAct pattern for multi-step reasoning and tool orchestration
- Deployment: Execute actions via the Tools layer with appropriate guardrails

Domain agents implement end-to-end flows documented in views based on the **End-to-End Scenario viewpoint** (Section 3.6), executing sequences that span data retrieval, decision-making, and action execution. They maintain domain-specific context while delegating infrastructure concerns to shared layers.

**Solution examples:** LangChain, Semantic Kernel (agent frameworks with tool binding and memory); GPT-4, Claude, Gemini Pro (LLM reasoning engines with function calling).

### 5.1.3 Layer 3: Tools (reusable actions)

The Tools layer provides reusable capabilities invoked by domain agents: API connectors, data retrieval functions, graph queries, validation logic, and action executors. This layer implements infrastructure decisions documented in views based on the **Physical/Infrastructure viewpoint** (Section 3.5), delivering 40-60% reusability across use cases by centralizing integration patterns.

Backend integration targets:

- **Knowledge Graph:** GOLD-tier structured domain knowledge (network topology, service specifications, configuration standards)
- SDN Platform: Network programmability and orchestration APIs
- Assurance Systems: Real-time telemetry, performance metrics, and fault data

**Solution examples:** Neo4j (knowledge graph), Milvus/Qdrant (vector search); NetworkX (graph algorithms for path computation); Prometheus, Elastic Stack (telemetry and observability); **MCP servers** (Model Context Protocol - standardized tool exposure for LLM integration with external systems).





### 5.1.4 Architectural principles

**Separation of concerns:** Orchestration logic remains independent from domain reasoning, which stays isolated from backend integration. This allows parallel development: platform teams enhance the Tools layer while domain experts refine agent logic.

**Framework alignment:** Each layer maps to specific framework constructs. The Orchestrator enforces governance and scoping boundaries; Domain Agents execute the four-phase lifecycle; Tools implement physical infrastructure capabilities. Views document these architectural decisions, ensuring choices in one layer remain traceable to artifacts captured in other viewpoints.

**Reusability mechanics:** Tools built for one use case (e.g., topology graph traversal for troubleshooting) become available to others (e.g., service design validation). By the fifth use case, 70-80% of the Tools layer requires no modification, reducing development time from weeks to days.

**Security model:** JWT-based authentication flows through all layers. Claims-based authorization ensures domain agents access only permitted backend systems. The Orchestrator validates tokens at ingress; Tools enforce backend-specific authorization at egress. This bidirectional security model implements audit requirements documented in governance views based on the Governance, Risk & Compliance viewpoint.

Practitioner's Perspective: Technology Selection and Three-Layer Benefits

Technology patterns examples observed across telecom deployments:

- Orchestration: LangGraph, custom FastAPI with intent classification
- **Agent frameworks:** LangChain or Semantic Kernel (function calling, memory, tool binding)
- LLM providers: OpenAI GPT, Anthropic Claude, Google Gemini
- **Knowledge layer:** Neo4j (graph topology), Milvus/Qdrant (vector search), NetworkX (path algorithms)
- **Integration standards:** MCP (Model Context Protocol) servers for standardized tool exposure, REST/gRPC for legacy systems
- **Observability:** Prometheus + Grafana (metrics), Elastic Stack (logs), OpenTelemetry (distributed tracing)

Why three layers? Building monolithic agents handling routing, domain logic, and integration initially appears faster, but multi-use-case deployments reveal systemic issues: agent competition on ambiguous requests, duplicated API integration code, and slow iteration cycles (routing changes require full redeployment). The three-layer architecture eliminates these bottlenecks - Orchestrator resolves ambiguity, shared Tools remove duplication, layer independence enables continuous deployment of domain logic without platform disruption. Tool integration evolution: Start with direct REST/gRPC connectors for rapid prototyping.





5.2. Use case validation 41

Migrate high-reuse integrations (Knowledge Graph, SDN APIs, telemetry systems) to MCP servers as the portfolio matures - standardized schemas reduce integration effort for new use cases and improve LLM tool discovery.

### 5.2 Use case validation

Three implementations demonstrate the architecture's versatility across the telecom operational spectrum: runtime incident response, design-time planning, and continuous optimization. Each use case shows how the four-phase lifecycle and architectural viewpoints guide concrete implementation decisions, with teams creating specific views that document their architectural choices within the relevant viewpoint frameworks.

### 5.2.1 Use case 1: Service recovery

**Problem:** Network path failures require rapid alternative routing while respecting operational constraints (e.g., avoid specific locations, upgrade bandwidth tiers).

### Framework mapping:

- Primary phase: Deployment (Section 2.5) runtime operations under SLA pressure
- **Key viewpoints:** Runtime/Process (incident response flows), Scenario (end-to-end path reconfiguration)
- Data tiers: GOLD (network topology graph), SILVER (real-time link status telemetry)

### **Architecture instantiation:**

- 1. User requests via interface: "Find new path avoiding Paris, upgrade to 10G"
- 2. Orchestrator authenticates request, routes to Network Operations domain agent
- 3. Domain agent executes multi-step plan:
  - Tool invocation: Query Knowledge Graph for topology structure
  - Tool invocation: Check Assurance Systems for current link health and KPIs
  - Internal reasoning: Compute feasible alternative paths using constraint solver
  - Tool invocation: Validate bandwidth availability against SDN platform inventory
  - Tool invocation: Execute path reconfiguration via SDN APIs
- 4. Domain agent returns confirmation with new path details and estimated restoration time

### Success KPIs (from Chapter 4):

- MTTR reduction
- Integration cost reduction: Single REST API vs. three separate operator portals
- SLA compliance increase







5.2. Use case validation 42

**Reusable components:** Knowledge Graph query tool (used in subsequent use cases), SDN connector (used in provisioning and optimization), path computation logic (adapted for design validation).

### 5.2.2 Use case 2: Network design assistant

**Problem:** Non-expert users struggle to create valid network topologies; design errors cause provisioning failures and project delays.

### Framework mapping:

- **Primary phases:** Scoping (requirements gathering) + Data (design validation against service specifications)
- **Key viewpoints:** Scenario (conversational design workflow), Governance (compliance with service rules and architectural standards)
- **Data tiers:** GOLD (service specifications, reference architectures), BRONZE (unstructured user requirements from meeting transcripts)

#### **Architecture instantiation:**

- 1. Conversational UI (Next.js + React Flow) collects requirements through LLM-guided dialogue
- 2. Domain agent validates design constraints using NetworkX-based graph rules engine
- 3. Tools layer retrieves applicable service specifications and reference architectures
- 4. React Flow visualization tool generates interactive topology diagram
- 5. AI proposes optimizations, prepares automated order submission to NaaS provisioning APIs

#### Success KPIs:

- Onboarding time for non-expert users reduction
- Design error rate reduction
- Proposal win rate improvement

**Reusable components:** Service specification validator (adapted for troubleshooting constraint checking), NetworkX-based graph rules engine (used in capacity planning), React Flow topology visualization (integrated into operational dashboards).

### 5.2.3 Use case 3: Automated troubleshooting

**Problem:** Manual troubleshooting requires scarce skilled operators, delays recovery, and increases operational costs.

### Framework mapping:

• **Primary phases:** Data (log and telemetry analysis) + Modeling (anomaly detection, root cause reasoning)







5.2. Use case validation 43

• **Key viewpoints:** Data/Model Lifecycle (telemetry pipelines and quality gates), Runtime (incident response orchestration)

• Data tiers: GOLD (device inventory and configuration baselines), SILVER (structured telemetry streams), BRONZE (unstructured error logs)

#### **Architecture instantiation:**

- 1. LLM-based agent interprets device status messages, identifies configuration drift and operational faults
- 2. Anomaly detection tool flags statistical deviations across device populations
- 3. Tools layer correlates logs, telemetry streams, and topology data to isolate fault domain
- 4. Network path analysis generates Tier-1 incident tickets with recommended remediation actions

#### Success KPIs:

- Average troubleshooting time reduction
- SLA compliance increase
- Tier-2 escalation rate reduction (more issues resolved at Tier-1)

**Reusable components:** Log parsing and normalization tool (used in security monitoring), telemetry correlation engine (adapted for performance optimization), ticket integration API (shared across operational workflows).

### 5.2.4 Cross-use-case patterns

**Table 5.1:** Framework adaptation across use cases

Element	Service Recovery	Network Design	Troubleshooting
Lifecycle	Deployment	Scoping + Data	Data + Modeling
emphasis			
Timeline	Minutes (incident)	Days (design cycle)	Minutes (MTTR)
Data	Topology, link status	Service specs,	Telemetry, logs,
sources		requirements	inventory
Domain	Path solver, SDN APIs	Graph validator, LLM,	Anomaly detector, LLM,
tools		visualizer	correlator
<b>GOLD</b> tier	Topology graph	Service specifications	Device inventory
data			
Success	MTTR	Onboarding time	Troubleshooting time
metric			

All three use cases share infrastructure skeleton despite serving different operational domains:





- Authentication: JWT with role-based claims and scope validation
- LLM gateway: Centralized prompt management, response caching, and cost tracking
- Monitoring: Unified observability for agent behavior, tool latency, and cost attribution
- Knowledge Graph: Common semantic layer providing domain context across use cases

This convergence validates the framework's reusability hypothesis: investing in solid shared infrastructure during the first two use cases delivers significant returns by use cases five through ten.

### 5.3 Implementation roadmap

Organizations adopting this framework can progress through four stages, balancing rapid value delivery with architectural sustainability. Timelines and resource allocation vary based on organizational readiness, domain complexity, and existing infrastructure maturity.

### 5.3.1 Stage 1: Pilot use case

Focus: Validate framework applicability through one well-scoped implementation.

- Choose high-value, bounded problem (service recovery, network design, troubleshooting)
- Build foundational three-layer architecture with minimal viable features
- Explicitly map implementation to framework: document phase transitions (Chapter 2) and populate initial views for primary viewpoints (Chapter 3)
- Establish baseline → target KPIs with measurement infrastructure (Chapter 4)

**Success criteria:** Agent successfully handles real user requests end-to-end; measurable improvement in target KPI; architectural patterns documented for replication.

### 5.3.2 Stage 2: Architecture refinement

Focus: Extract reusable patterns from pilot experience.

- Identify common components suitable for platform layer (authentication, logging, error handling)
- Document domain agent template and tool invocation patterns
- Establish governance foundations: security model, audit requirements, compliance checkpoints
- Create capability roadmap: build-vs-buy decisions for anticipated use cases

**Success criteria:** Second use case development time significantly reduced through component reuse; governance controls operational.







5.4. Conclusions 45

### 5.3.3 Stage 3: Use case expansion

**Focus:** Apply validated architecture to additional domains.

• Deploy use cases following established patterns (network planning, customer analytics, security operations)

- Track infrastructure reuse metrics to validate platform maturity
- Gate new use cases on demonstrable KPI improvement or documented learnings
- Establish formal change management as portfolio grows (impact analysis, architecture review, incident retrospectives)

**Success criteria:** Infrastructure reuse increases with each deployment; new use cases launch faster than pilot; platform changes don't destabilize existing use cases.

### 5.3.4 Stage 4: Platform maturation

**Focus:** Transition from collection of use cases to unified agent platform.

- Consolidate similar tools and connectors based on observed patterns
- Enhance orchestrator intelligence (intent classification, context management, cross-domain routing)
- Optimize operational capabilities (monitoring, cost controls, performance tuning)
- · Evolve governance policies based on production incidents and audit feedback

**Maturity indicator:** New use cases deploy rapidly with minimal platform modifications. Focus shifts from proving feasibility to optimizing efficiency.

### 5.4 Conclusions

The three-layer architecture provides a proven pattern for translating framework concepts into operational systems:

- Orchestrator enforces governance boundaries and routing logic across all use cases
- Domain Agents execute the four-phase lifecycle within specialized domains
- Tools centralize integrations, enabling component reuse across implementations

Architectural views document decisions at each layer: governance views capture policy enforcement in the Orchestrator, scenario views trace execution flows through Domain Agents, and infrastructure views specify tool implementations and backend integrations. This documentation ensures architectural choices remain visible, traceable, and revisable as the system evolves.

Use case validation across service recovery, network design, and troubleshooting demonstrates the framework's versatility: consistent structure applied to diverse domains yields measurable outcomes. Organizations that master this architecture on their first use case accelerate subsequent deployments through disciplined component reuse and pattern recognition.





5.4. Conclusions 46

The framework's ultimate value lies not in any single implementation but in establishing architectural discipline that scales: governance mechanisms that prevent shortcuts, data practices that ensure quality, and operational patterns that compound returns across use cases.





# 6 Summary

This framework exists to accelerate telecom AI adoption by avoiding common pitfalls. The preceding chapters presented a pragmatic solution: a four-phase execution lifecycle integrated with six architectural viewpoints, validated through real-world use case deployments. This chapter validates the framework against the obstacles identified in Chapter 1 and synthesizes lessons for organizational adoption.

## 6.1 Framework validation: problems solved

The Introduction (Chapter 1) identified core failure modes that trap telecom AI initiatives in proof-of-concept purgatory. This section validates how the framework addresses each category of challenge: delivery execution obstacles first, then the governance and risk controls that enable enterprise adoption.

### 6.1.1 Execution challenges: delivery and integration

The following table maps core execution failures to framework solutions that restore predictability and measurable progress.

Table 6.1: Execution validation: Delivery and integration obstacles addressed through framework design

Challenge	Framework solution	Described examples
Undefined ROI - Pilots lack baseline KPIs and	KPI Baseline + Targets (Section 2.2, Chapter 4)	Phase 1 Scoping establishes measurable success criteria before building
target improvements		
Data Overreach - "Knowledge lake" initiatives delay impact	GOLD Seed Curation (Section 2.3, Section 3.1)	Start with a small GOLD set; Data phase emphasizes quality over volume





Challenge	Framework solution	Described examples
Fragile Prototypes - Notebooks and scripts create operational fragility	Evaluation Harness + Architecture (Section 2.4, Section 3.3, Section 5.1)	Automated quality gates block regressions; version control for prompts, tools, and models; three-layer architecture (Orchestrator → Domain Agents → Tools) provides production-grade separation of concerns
Disconnected ML Assets - Existing models remain siloed from agent workflows	ML Tool Integration (Section 2.4, Section 5.1.3)	ML forecasting, classification, anomaly detection exposed as callable tools with strict JSON schemas; Tools layer (Section 5.1.3) provides standardized integration pattern; troubleshooting use case demonstrates anomaly detector integration for automated ticket generation

These execution patterns - starting narrow with clear KPIs, curating quality over volume, automating quality gates, integrating existing ML assets - restore predictability to AI delivery. However, solving execution challenges alone doesn't unlock enterprise production deployment. In telecom operations, governance and risk controls determine whether pilots scale or stall.

### 6.1.2 Governance & risk controls: the enterprise enabler

Governance isn't an afterthought - it's the difference between production approval and permanent pilot status. Enterprise AI governance embeds risk controls throughout the lifecycle. The framework addresses governance as a distinct challenge category because telecom operators face regulatory scrutiny, compliance mandates, and operational risk thresholds that consumer-focused AI deployments rarely encounter.

The following table details how each enterprise risk category maps to framework mitigation strategies, expanding on the governance challenge introduced above.





 Table 6.2: Governance validation: Risk mitigation through architectural implementation

Risk category	Mitigation approach	Described patterns
Hallucination & Misinformation	Citation requirements, validator agents, evaluation gates	Quality gates in Modeling phase (Section 2.4) require cited sources; governance views based on the Governance, Risk & Compliance viewpoint (Section 3.2) document traceability requirements; Multi-agent architecture (Section 5.1) includes validation layer with groundedness scoring
Data Leakage &	Tiered data access, PII	Data & Model Lifecycle views (Section 3.1)
Privacy	redaction, audit logs	document GOLD/SILVER/BRONZE tiers with role-based permissions; immutable trace logs documented in Runtime/Process views (Section 3.3); Phase 1 risks (Section 2.2.3) enforce encryption and redaction at API layer
Cost Inflation	Cost budgets, token limits, model tiering	Deployment phase (Section 2.5) monitors cost per 1K requests; Financial KPIs (Chapter 4) track telecom-specific savings (truck roll avoidance, AHT reduction); Architecture (Section 5.1) enables fast/slow model selection based on request complexity
Vendor Lock-In	Abstraction layers, standardized tool schemas	Multi-agent architecture (Section 5.1.3) demonstrates gateway abstraction enabling provider swaps; Phase 1 risks (Section 2.2.3) emphasize abstraction layers; MCP (Model Context Protocol) servers standardize tool interfaces across LLM providers
Model/Data Drift	Scheduled evaluation, drift detection, automated alerts	Data & Model Lifecycle views (Section 3.1) document drift thresholds; Quality KPIs (Section 4.2) monitor feature/embedding drift; automated review triggers quarantine stale data; separate drift layers for ML features vs. text embeddings
Compliance & Audit Gaps	Immutable trace logs linking full decision chain	End-to-End Scenario views (Section 3.6) map every request to versioned decisions documented across all viewpoints; governance views (Section 3.2) document audit requirements with trace IDs, prompt hashes, model versions, and policy enforcement timestamps





Counter-intuitively, governance accelerates adoption rather than slowing it. Stakeholders trust systems they can audit: finance approves budgets with cost attribution, compliance unblocks production deployment when trace logs exist.

The cumulative impact of addressing both execution and governance challenges - through explicit KPIs, curated knowledge, automated quality gates, governed evolution, and embedded risk controls - transforms the framework from a technical guide into an organizational enabler. Organizations conclude AI is immature when structural discipline is missing, not when technology falls short.

### 6.2 Final recommendations

The framework validation above demonstrates how intentional design addresses telecom AI adoption obstacles. Translating these patterns into organizational action requires strategic focus. The following recommendations prioritize the choices that most consistently differentiate successful deployments from stalled pilots.

### Start narrow, prove value fast:

Pick one high-value process, define 3-5 measurable KPIs, curate a focused GOLD knowledge set (100-500 documents), and prove measurable uplift in 60-90 days. The flywheel (Section 1.2) begins with a scoped KPI gap - not a strategic vision, but a specific, testable hypothesis about where AI can improve a concrete operational process.

### Governance enables speed:

Embedding controls from day one accelerates adoption. As demonstrated in the governance validation above (Section 6.1.2), controls transform from obstacles into enablers when architected intentionally.

### Design for reusability:

Invest in architectural patterns that enable multi-use-case reuse from the start. The three-layer pattern (Section 5.1) - Orchestrator  $\rightarrow$  Domain Agents  $\rightarrow$  Tools - demonstrates this principle: shared infrastructure components (authentication, LLM gateway, monitoring, knowledge retrieval) remain stable across use cases while domain-specific logic evolves independently. This separation of concerns enables faster second and third deployments by amortizing infrastructure investment across multiple operational processes.

### **Evolve views incrementally:**

Start with three core viewpoints (Data & Model Lifecycle, Runtime, Scenario), populate initial views for each, add Governance and Physical views pre-deployment, and layer Capability views as operational patterns emerge (Section 3.7). Let production telemetry guide evolution. Views are living documents - two pages at pilot launch, ten pages six months into production, enriched with operational learnings and incident post-mortems.







### 6.3 Closing perspective

The framework delivers on its core promise: breaking the telecom AI adoption trap through intentional design. Not by waiting for perfect technology or comprehensive planning, but by starting narrow, measuring rigorously, governing deliberately, and compounding value through each iteration. The difference between stalled pilots and production success isn't technical maturity - it's structural discipline.

Organizations that adopt this framework systematically won't just deploy agents - they'll build organizational capability to evolve AI systems safely, measurably, and sustainably. That capability, more than any single implementation, represents the framework's ultimate value.





# A ML Workflow Blueprint - Key Aspects

This blueprint outlines a practical, step-by-step workflow for designing ML systems. It is based on the "Machine Learning in Production" course by Andrew Ng, which offers a useful conceptual baseline.

### A.1 Phases descriptions

### A.1.1 Scoping (define project)

- Decide which problem to solve with ML.
- Define X (inputs) and Y (outputs).
- Establish clear success criteria what does "good enough" look like?

### A.1.2 Data (define, label, organize)

- Define data requirements and assemble a baseline dataset.
- Collect, clean, and label data this is often the most time-consuming stage.
- Ensure consistent formats, accurate labels, and representativeness.

### A.1.3 Modeling (select, train, error analysis)

- Select appropriate ML model(s).
- Train iteratively and document experiments.
- Perform error analysis to identify systematic failures or bias.
- Loop back to data collection or model changes as lessons emerge.

### A.1.4 Deployment (production, monitoring)

- Deploy the system into production and integrate with applications or services.
- Monitor live performance, data drift, and failures.
- Maintain the system: update models when distributions shift and retrain as needed.
- Treat deployment as a midpoint rather than an endpoint long-term maintenance is critical.







### A.2 Do's and don'ts by stage

### A.2.1 Scoping

#### Do

- Clearly define the problem and success metrics.
- Confirm ML is the right tool some tasks are better solved with rules or heuristics.
- Keep scope realistic.

#### Don't

- Jump into modeling without clarity on the problem.
- Set vague goals (for example, "improve accuracy") without measurable criteria.

### A.2.2 Data

#### Do

- Collect diverse, representative data.
- Label consistently and audit annotations for errors.
- Split datasets properly (train/validation/test).
- Establish a baseline prior to training.

### Don't

- Assume more data always yields better results.
- Ignore class imbalance or rare categories.
- Mix training and test data (avoid leakage).

### A.2.3 Modeling

#### Do

- Start with simple models before increasing complexity.
- Use error analysis to guide improvements and spot systematic mistakes.
- Iterate between model design and data fixes.
- · Record experiments and hyperparameters for reproducibility.

#### Don't

- Overfit by optimizing only for validation accuracy.
- Treat modeling as a one-off activity.
- Ignore interpretability when stakeholders require it.





### A.2.4 Deployment

#### Do

- Monitor for data drift and model degradation.
- Automate retraining and updates where practical.
- Log predictions and outcomes to close feedback loops.
- Treat maintenance as an ongoing responsibility.

### Don't

- Assume deployment equals completion real-world feedback will change behavior.
- Overlook scalability, latency, or cost constraints.
- Neglect user feedback or ethical concerns such as fairness and privacy.



