## Al Agents: From Concept to Real-World Impact

Nitin Kumar — Director, Data Science (GenAl)



# "Everyone's talking about agents."

Boards. Blogs. Hallways. Demos.

"Let's separate hype from impact."

### But here's the truth...

# Most organizations are still stuck in POCs and pilots.

We'll talk about how to move to production—safely and repeatedly.

### What is an Al agent?

**Definition:** A goal-directed system that *perceives* inputs, *reasons* with context, and *acts* via tools to achieve outcomes.

Core components: LLM brain • Tools/Environment • Memory/State • Guardrails • Observability • Human-in-the-Loop.

#### **Tools & Environment**

Search & Retrieval (web, RAG, vector DBs)

APIs & Data Access (email/calendar/CRM; SQL/warehouse)

#### **Memory & State**

Short-term: conversation context (billed in LLM call)

Long-term & semantic: vector DB retrieval; preferences/rules

#### Guardrails

Policies & PII redaction
Schema-constrained outputs & tool gating

#### Observability

Traces & metrics (latency, cost, accuracy)
Audit logs & incident response

Human-in-the-Loop (HITL)

Shadow → assisted → autonomous Approvals & escalation paths

Say it once: **perceive**  $\rightarrow$  **reason**  $\rightarrow$  **act** ... powered by tools & memory, bounded by guardrails, visible via observability, and partnered with humans.

# The Two-Layer Framework — Layer 1: Architecture (The Body)

#### Single Agent

One agent handles the entire task end-to-end.

#### **Sequential Pipeline**

Assembly-line stages with clear handoffs.

#### **Parallel Multi-Agent**

Specialists work simultaneously; merge results.

#### **Orchestrator**

Manager delegates to specialized workers.

#### **Collaborative**

Peer agents collaborate without a hierarchy.

#### **Network / Graph**

Conditional routing across a graph of agents.

Pick the simplest structure that meets your risk, scale, and governance needs.

# The Two-Layer Framework — Layer 2: Reasoning (The Brain)

Zero-/Few-Shot

Direct responses with optional exemplars.

**Chain-of-Thought** 

Step-by-step reasoning before answering.

ReAct

Loop: Reason → Act → Observe for tool-use.

Plan-and-Execute

Plan upfront, then execute systematically.

Reflexion

Self-critique and iterative improvement.

**Tree-of-Thoughts** 

Explore multiple reasoning paths; pick best.

choose *one* from each layer — and you've defined your agent's behavioral pattern.

## Four Critical Challenges to Overcome

#### **Data Quality**

Accuracy • Completeness • Timeliness

Siloed data, legacy formats

Inconsistent standards

Missing business context

#### **Hallucinations & Reliability**

Grounding • Guardrails • SLAs

Made-up facts & brittle prompts

Mitigate with RAG, constraints, evals

Measure uptime & accuracy SLOs

#### **Compliance & Governance**

Operate within guardrails

Privacy (GDPR/CCPA), security, audit trails

Bias detection & transparency

Use-case tiering, model cards

#### **Trust & Adoption**

People • Process • Change

HITL workflows & clear escalation

Training, comms, incentives

Shadow → assisted → autonomous

Visual: think of these as four pillars (or puzzle pieces) that lock together.

## Garbage In, Garbage Out Still Applies

The data quality triad

**Accuracy** 

Completeness

**Timeliness** 

**Common issues** 

Siloed data

Legacy formats

Inconsistent standards

Missing context

**Solution approach** 

Data foundations before Al ambitions Authoritative sources, contracts, lineage

Freshness SLOs, metadata, stewardship

Privacy-by-design (masking, access)

fund data contracts & lineage before model spend.

## When Al Gets Creative (And That's Bad)

#### What & why

Made-up facts or unsupported claims

Causes: weak grounding, vague prompts,
adversarial inputs

#### Mitigation strategies

Retrieval-Augmented Generation (RAG) with citations

Confidence scoring & abstain/fallback

Human-in-the-loop validation for edge cases

Constrained outputs (tools, schemas)

Testing & monitoring in production

require sources for claims; measure answerability not just accuracy.

## Operating Within Guardrails

Regulatory landscape

GDPR • CCPA • industry regs

**Key considerations** 

Data privacy

Explainability / transparency

Audit trails

Bias detection

Security

**Framework** 

Responsible Al principles in practice:

Purpose limitation & data minimization

Role-based access & regional routing

Policy checks & model cards

Incident response & red-teaming

partner with Legal & Privacy early; document decisions & audits.

## **Building Trust Through Consistency**

The reliability stack

Robust testing frameworks (gold sets, adversarial)

Monitoring & observability (trace, cost, latency, drift)

Graceful degradation (fallbacks, safe defaults)

Feedback loops (HITL review → retrain)

Version control & rollback

Ask the hard question

What's your agent's uptime & accuracy SLA?

Define SLOs, error budgets, and auto-rollback gates

ship monitors before launch; tie alerts to rollbacks.



## Augmentation, not replacement.

Speed & scale of Al × empathy & judgment of humans.

humans remain accountable; design clear escalation paths.

## Human-in-the-Loop by design

Tier 0

Autonomous for low-risk, high-volume tasks with fallbacks.

Tier 1

Suggest & draft; humans approve, edit, or escalate.

Tier 2

Expert-only: Al assists with retrieval, summarization, rationale.

Escalation paths, audit trails, and feedback loops are non-negotiable.

## What We've Learned: Start Small, Think Big

**Start with clear, measurable problems** — not technology looking for problems.

**Build trust incrementally** — shadow mode → assisted mode → autonomous.

**Invest in data infrastructure first** — foundation before flashiness.

start tiny, prove value, then earn autonomy.

## What We've Learned: People & Process Matter

Change management is not optional — train, communicate, iterate.

**Design for observability** — you can't improve what you can't measure.

**Plan for the unexpected** — edge cases will find you.

incentives drive adoption; metrics drive improvement.

## **Metrics That Matter**

**Technical** 

Accuracy

Latency

**Uptime** 

**Business** 

ROI

Time saved

Volume handled

User

Adoption rate

Satisfaction

Trust scores

Use a balanced scorecard across technical, business, and user lenses.

## From Talk to Action: Your Next Steps

#### Starting out

Identify one high-value, low-risk use case

Assemble a cross-functional team

**Build data foundations** 

Start with an augmentation mindset

#### Scaling up

Standardize frameworks & patterns

Build centers of excellence

Share learnings across teams

assign an owner, a metric, and a deadline—today.

## Remember These Three Things

**Reality over hype** — focus on deployment, not just experimentation.

**Augmentation mindset** — Al + Human = best outcomes.

**Foundations matter** — data quality, reliability, and trust are non-negotiable.

three anchors—Reality. Augmentation. Foundations.

## Thank you

Connect with me on LinkedIn



Tap • Click • Arrow keys