

# GTC 2026 — Article B5

## Enterprise AI Adoption Guide

GTC 2026 Series · BOTUM Analysis · March 2026

March 2026

## Table of Contents

---

1. Introduction — From interest to production
2. Phase 1: Evaluation — Identifying the right use cases
3. Phase 2: Pilot — Validating with a concrete project
4. Phase 3: Scale — Going from 1 to N agents in production
5. Phase 4: Governance — Control, audit, secure
6. Classic mistakes and how to avoid them
7. Final checklist — 20 questions before go-live

## Introduction — From interest to production

At GTC 2026, NVIDIA confirmed that agentic AI is no longer a lab topic: it is a production reality. NVIDIA Agent Intelligence Toolkit, NIM, NemoClaw, agentic Blueprints — the building blocks are here. The real question is: how does your organization move from curiosity to deployment?

This guide targets solution architects and CIOs who need to make concrete decisions. It structures the journey into 4 phases: Evaluation, Pilot, Scale, Governance. Each phase includes precise actions, success criteria, and pitfalls to avoid.

i BOTUM note: this guide is based on GTC 2026 announcements and our experience guiding Canadian enterprises through AI projects. Figures cited are order-of-magnitude estimates valid for standard enterprise contexts.

## Phase 1: Evaluation — Identifying the right use cases

Most enterprise AI projects fail not because of technology, but because of a poor use case selection at the start. Evaluation must be methodical.

### Pilot use case selection criteria

Criteria	Good sign	Bad sign
Task volume	> 100 repetitive tasks/day	Rare or highly variable tasks
Data quality	Structured data available	Siloed or missing data
Error tolerance	Error = negligible cost	Error = legal or security risk
Measurability	Clear KPI (time, cost, quality)	Diffuse, non-measurable impact
Business support	Business champion identified	IT project without business sponsor
AI complexity	Simple orchestration (2-3 agents)	Multi-agent critical from day 1

### Valid enterprise use cases from GTC 2026

- Tier-1 customer support: NIM agent + RAG on internal knowledge base (ROI 3-6 months)
- Document analysis: contracts, RFPs, reports — Blueprint agent with local NIM LLM
- Proactive IT monitoring: monitoring agent + automatic escalation via NeMo Guardrails
- Employee onboarding: multi-step guided journey with human validation checkpoints
- Report generation: data aggregation + automated drafting with human supervision

i BOTUM signal: enterprises that succeed choose use cases where the agent augments a human rather than fully replacing them. Internal adoption is 3x faster.

### Evaluation phase deliverables (4-6 weeks)

- Vision document: 1 page, selected use case, target KPIs, initial scope

- Data audit: source inventory, quality, accessibility, GDPR compliance
- Provisional target architecture: NIM + orchestrator + interface stack
- Pilot budget: cloud GPU costs, development, integration estimates
- Validated sponsor: CIO or business director committed to the 90-day pilot

## Phase 2: Pilot — Validating with a concrete project

A pilot is not a POC. A POC proves the technology works. A pilot proves that your organization can operate it. The difference is fundamental.

### Recommended technical stack for the pilot

Component	BOTUM Recommendation	Alternative
LLM inference	NVIDIA NIM (cloud or on-prem)	OpenAI / Anthropic API
Agent orchestration	LangGraph + NVIDIA Agent Toolkit	CrewAI, AutoGen
RAG / memory	FAISS or Milvus + NeMo Retriever	Chroma, Weaviate
Guardrails	NeMo Guardrails (mandatory)	Regex + custom filter (insufficient)
Observability	Langfuse or LangSmith	Custom JSON logs
Interface	Gradio or Streamlit (pilot)	Existing business interface
Infrastructure	AWS or Azure (on-demand GPU)	On-prem if > 50k req/day

### Pilot success criteria (90 days)

- Performance: task success rate > 80% (defined before launch)
- Adoption: > 70% of target users actively using it after 30 days
- Costs: cost per agent task < cost per human task (or freed time > 20%)
- Reliability: uptime > 99% over the last 30 days of the pilot
- Guardrails: 0 incidents of inappropriate content or data leakage

i If the pilot does not meet these thresholds, do not move to Scale phase. Return to Evaluation with your learnings.

## Phase 3: Scale — Going from 1 to N agents in production

The transition from pilot to production is where most enterprise AI projects stall. The technology works, but the organization is not ready. Here is how to clear that hurdle.

### Production architecture: what changes vs. the pilot

- High availability: 2 regions minimum, automatic failover, 99.9% SLA

- GPU autoscaling: ability to scale 10x in < 5 minutes (cloud mandatory)
- Agent CI/CD pipeline: zero-downtime deployment, automated tests, instant rollback
- Production observability: distributed traces, latency/error alerts, business dashboards
- Semantic cache: reuse similar responses = -40 to -60% GPU costs
- Rate limiting: per-user/service quotas to protect the infrastructure

### Organizational model: the AI Ops team

Role	Responsibility	Profile
AI Ops Lead	SLA, incidents, GPU budget	Senior DevOps + LLM training
Prompt Engineer	Prompt optimization, evals	Python dev + linguistics
Data Steward	Data quality, GDPR, RAG	Data analyst + legal
Business Owner	Business KPIs, prioritization	Business director or manager
AI Security	Guardrails, audit, red team	SecOps + adversarial AI training

### Cost management at scale

- NIM microservices: deploy only the models in use (avoid one-size-fits-all)
- Model routing: small models for simple tasks, large models for complex ones
- Batching: group non-real-time requests to maximize GPU utilization
- Reserved instances: if stable volume > 60% of the time, cloud reservations = -40% cost

i BOTUM benchmark: a well-optimized customer support agent handles 500-800 conversations/hour on a single H100 GPU. Calculate your break-even vs. human agents at that ratio.

## Phase 4: Governance — Control, audit, secure

Governance is not a phase that comes after production: it must be embedded from the pilot. But it is in production where it becomes critical.

### BOTUM enterprise AI governance framework

- Acceptable use policy: what agents can and cannot do (documented, legally validated)
- Agent registry: inventory of all agents in production, their access, their capabilities
- Complete audit trail: every agent decision traced, timestamped, exportable for audit
- Human-in-the-loop: human validation process for high-impact decisions
- Regular red teaming: quarterly adversarial prompting tests, report to leadership
- Model update process: validation procedure before every LLM update in production

### GDPR and data protection compliance

- Never send personally identifiable information (PII) to an external LLM without consent
- Local NIM or private VPC for sensitive data: healthcare, finance, HR
- Right to explanation: if an agent makes a decision affecting an individual, explanation is mandatory
- Log retention: define agent trace retention period (recommended: 12 months)
- DPO involvement: every new agent capability must go through a Data Protection review

## Classic mistakes and how to avoid them

Mistake	Symptom	Solution
Too big too fast	Pilot on 20 use cases in parallel	Max 2 use cases in phase 1
No guardrails	Agent answers anything out of scope	NeMo Guardrails mandatory
Unprepared data	RAG hallucinating 30% of the time	Data audit before pilot
No business champion	Adoption < 20% after 3 months	Business owner required
Model lock-in	Total dependency on GPT-4	LLM abstraction (LiteLLM)
Ignoring GPU costs	Cloud bill 5x vs estimate	GPU FinOps from pilot day 1
No observability	Impossible to debug in prod	Traces + dashboards from day 1
Forgetting security	Prompt injection by users	Adversarial tests pre-launch

## Final checklist — 20 questions before go-live

### Evaluation & Architecture

- Does the use case have a measurable KPI and an identified business sponsor?
- Is the data audit complete (quality, scope, GDPR)?
- Is the NIM + orchestrator + RAG architecture documented and validated?
- Have GPU costs been estimated and a cloud budget approved?
- Is a human fallback in place for edge cases?

### Security & Compliance

- Is NeMo Guardrails or equivalent configured and tested?
- Have prompt injection tests been conducted?
- Is the PII processing policy documented and validated by the DPO?
- Is the audit trail active and exportable?
- Is the agent registry created and maintained?

## Operations & Scale

- Is observability (traces, metrics, alerts) in place?
- Is a CI/CD pipeline with automated tests configured?
- Is the incident runbook (what to do when the agent fails) written?
- Is the AI Ops team trained and on-call rotations defined?
- Has the rollback procedure been tested?

## Adoption & Governance

- Is end-user training planned?
- Is a user feedback process in place?
- Has the acceptable use policy been communicated to all?
- Is a quarterly AI review committee scheduled?
- Is the scale-up plan (Scale phase) documented?

BOTUM supports your enterprise AI adoption

From use case evaluation to production deployment, BOTUM teams guide CIOs and solution architects at every step. Free audit available for Canadian organizations. Contact:

[www.botum.ca/contact](http://www.botum.ca/contact)

-> Online version: [blog.botum.ca/gtc2026-b5-enterprise-ai-adoption-guide/](https://blog.botum.ca/gtc2026-b5-enterprise-ai-adoption-guide/)