

# NexaStack Solution for Azure Marketplace

Problem Statement .....	2
Solution Detail .....	2
Key Components .....	4
Integration Points.....	5
Use Cases .....	5
Customer Pain Points Addressed .....	5
Industry-Specific Applications .....	6
Sample Customer Journey .....	7
Technical Requirements .....	7
Performance Considerations .....	8
Tools and Azure Services Used .....	9
Users of Agent .....	9
Dependencies .....	10
Key Benefits and Differentiators .....	11
Value Proposition.....	12
Conclusion .....	12

NexaStack is a unified AI infrastructure platform that simplifies how enterprises deploy, manage, and scale models across environments—cloud, on-premises, and edge. It supports LLMs, vision, and multimodal AI while offering seamless orchestration under one stack. By unifying diverse workloads, NexaStack ensures agility, performance, and efficiency at scale.

It helps to accelerate enterprise AI adoption by providing a consistent and reliable infrastructure for running production AI workloads. It enables deployment of any model, integrates observability for monitoring performance, and offers cost optimization features to maximize ROI. Enterprises can scale AI faster while maintaining operational resilience and governance.

## Problem Statement

Enterprises adopting agentic AI face several operational challenges:

- **Unreliable Agent Behavior:** Agents may drift, hallucinate, or fail due to LLM unpredictability.
- **Lack of Observability:** Limited visibility into agent reasoning, memory usage, and tool interactions.
- **High Operational Risk:** No structured rollback, SLA tracking, or incident response mechanisms.
- **Security & Compliance Gaps:** Exposure to prompt injection, unauthorized API usage, and missing audit trails.
- **Scaling Difficulty:** Hard to manage multiple agents, coordinate workflows, and ensure governance.

Without an AI MSP framework, enterprises risk **downtime, compliance failures, security vulnerabilities, and uncontrolled costs**.

## Solution Detail

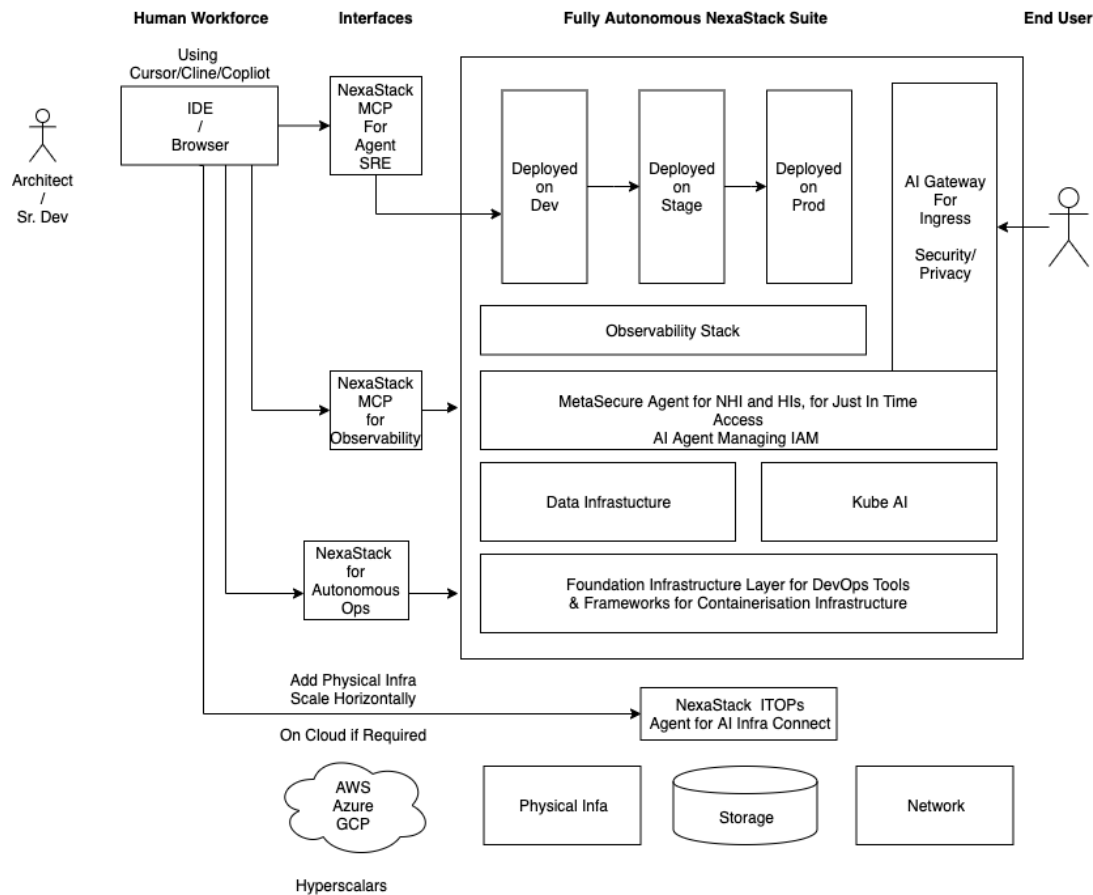
NexaStack addresses these challenges by providing a **holistic AI MSP practice** built around ten pillars:

1. **Agent Lifecycle Management** – deployment, versioning, safe rollback.

2. **Uptime & Reliability** – SLA enforcement, fallback mechanisms, failure handling.
3. **AI Observability & Tracing** – telemetry, cognitive traces, token/GPU cost monitoring.
4. **Drift & Anomaly Detection** – hallucination detection, tool misuse prevention.
5. **AI Security & Access Control** – sandboxing, ACLs, jailbreak protection.
6. **Compliance & Responsible AI (RAI)** – audit logging, bias monitoring, governance.
7. **Agent Coordination & Scheduling** – multi-agent orchestration with traceable inter-agent calls.
8. **Platform Operations** – ModelOps, GPU/infra management, memory ops.
9. **Feedback & Continuous Improvement** – human-in-the-loop (HIL), prompt versioning, regression testing.
10. **Customer Support & Incident Response** – runbooks, RCA templates, SLA-based support.

By aligning with these pillars, NexaStack delivers **reliable, secure, and observable AI operations at scale**.

## Technical Architecture



## Key Components

- **AgentOps Stack:** Covers lifecycle, uptime, rollback, and SLA management.
- **Observability Platform:** LangFuse + OpenTelemetry for agent traces, cost dashboards, and incident tracking.
- **Security Guardrails:** Sandbox execution, ACLs, prompt injection defense.
- **Anomaly Detection Engine:** Drift detection and hallucination monitoring using OPERA-based pipelines.
- **Orchestration Layer:** A2A + LangGraph workflows for multi-agent collaboration.
- **Compliance Module:** Automated audit logs, governance policies, and RAI monitoring.
- **Support Framework:** Playbooks, RCA workflows, and structured customer success enablement.

## Integration Points

NexaStack integrates seamlessly into enterprise IT and AI ecosystems:

- **LLMs & Models:** Mixtral, LLaMA, proprietary/custom models with autoscaling GPU support.
- **Data Systems:** ERP (SAP, Oracle), CRM (Salesforce), Data Lakes (Snowflake, Databricks).
- **Observability Tools:** LangFuse, OpenTelemetry, Prometheus, Grafana.
- **CI/CD Pipelines:** Integration with DevOps workflows for continuous agent updates.
- **Security & IAM:** Enterprise RBAC, identity providers (Okta, Azure AD).
- **Compliance Systems:** Audit logs aligned with enterprise GRC frameworks.

## Use Cases

NexaStack enables enterprises to adopt AI agents in mission-critical domains:

- **FinOps** – Cloud cost optimization, tagging hygiene, budget enforcement.
- **SRE** – Automated RCA, incident response, drift detection.
- **Compliance & Governance** – Continuous audit readiness, policy enforcement.
- **DevOps** – Self-optimizing CI/CD rollouts, rollback orchestration.
- **AI Governance** – Bias monitoring, safe model deployment, trace-based risk evaluation.
- **Customer Success** – AI-driven support workflows with SLA guarantees.
- **Industry-Specific:** Banking (regulatory audits), Insurance (claim automation), Manufacturing (predictive maintenance), Public Safety (real-time insights).

## Customer Pain Points Addressed

NexaStack directly tackles the most pressing challenges enterprises face when adopting AI agents:

- **Unreliable Performance:** Agents that drift, hallucinate, or fail without rollback safety nets.
- **Black-Box Operations:** Lack of visibility into agent decisions, reasoning, and cost attribution.
- **Operational Complexity:** Difficulty managing multi-agent workflows, tool integrations, and SLA compliance.
- **Security Vulnerabilities:** Exposure to prompt injection, unauthorized API/tool calls, or unmonitored execution.
- **Compliance Risks:** Absence of audit trails, bias monitoring, and regulatory readiness.
- **Cost Uncertainty:** Poor tracking of GPU usage, token consumption, and scaling costs.

By embedding **observability, reliability engineering, orchestration, and RAI guardrails**, NexaStack eliminates these bottlenecks and de-risks enterprise adoption.

## Industry-Specific Applications

NexaStack powers domain-tailored AI agent use cases across industries:

- **Banking & Financial Services**
  - AI agents for **regulatory compliance monitoring** (RBI, Basel norms).
  - Automated **FinOps** for cloud cost governance.
  - Intelligent **fraud detection & audit reporting**.
- **Insurance**
  - Claims automation and validation with audit-ready decision trails.
  - Compliance-driven risk assessments for underwriting.
  - Multi-agent workflows for customer service (policy queries, payments).
- **Manufacturing & Industrial Automation**
  - Predictive maintenance agents using telemetry + anomaly detection.
  - Supply chain optimization agents coordinating with ERP systems.
  - Safety compliance agents monitoring equipment and workforce.
- **Public Safety & Law Enforcement**
  - Real-time anomaly detection for surveillance and IoT feeds.
  - AI-driven audit trails for evidence management.
  - Multi-agent orchestration for incident coordination.

- **Healthcare & Life Sciences**
  - Clinical decision support agents with **RAI guardrails**.
  - Compliance monitoring for HIPAA/GDPR.
  - Predictive patient flow & scheduling agents.

## Sample Customer Journey

- 1. Discovery & Assessment**
  - a. NexaStack team evaluates customer systems (ERP, CRM, Data Lakes) and identifies candidate use cases.
  - b. Security & compliance requirements are mapped.
- 2. Pilot Engagement**
  - a. Deploy 1–2 agents in a controlled environment (e.g., FinOps, SRE).
  - b. Enable observability (LangFuse + OpenTelemetry) for trace-based validation.
  - c. Establish rollback & SLA monitoring baselines.
- 3. Scale-Up Phase**
  - a. Expand to 5+ workflows with orchestration (LangGraph, A2A).
  - b. Introduce anomaly detection & compliance monitoring modules.
  - c. Integrate with enterprise IAM and DevOps pipelines.
- 4. Enterprise Rollout**
  - a. Standardized onboarding via playbooks, checklists, and runbooks.
  - b. SLA-backed managed service contract (Startup, Growth, Enterprise tiers).
  - c. Continuous improvement loop with HIL feedback and regression testing.
- 5. Long-Term Partnership**
  - a. Cost dashboards and GPU usage optimization.
  - b. Automated compliance reports for audits.
  - c. Joint innovation for new AI-driven industry use cases.

## Technical Requirements

To operationalize NexaStack, the following technical prerequisites apply:

- **Infrastructure**
  - Cloud (AWS/Azure/GCP) or On-Premises GPU clusters.

- Supported GPUs: T4, A10G, L40S, A100, H100, H200 (scalable by workflow).
- Storage for memory ops (Redis, mem0, Vector DBs).
- **Software & Tooling**
  - Agent Orchestration: LangGraph, A2A, MCP.
  - Observability: LangFuse + OpenTelemetry + Prometheus.
  - Security: Sandbox execution, ACL/IAM integration, prompt injection protection.
  - Compliance: Audit logging, RAI classifiers, lineage tracking.
- **Enterprise Integrations**
  - ERP/CRM: SAP, Oracle, Salesforce, Odoo.
  - Data Platforms: Snowflake, Databricks, S3, Azure Data Lake.
  - CI/CD Pipelines: GitHub Actions, GitLab, Jenkins.
  - IAM & Security: Okta, Azure AD, enterprise PKI.
- **Operational Readiness**
  - SLA monitoring (uptime dashboards, latency tracking).
  - Rollback registry with previous stable agent states.
  - Multi-tenant support for internal teams or external customers.

## Performance Considerations

NexaStack is designed with enterprise-scale performance in mind, balancing **latency, reliability, and cost efficiency**:

- **Uptime Guarantees:** SLA-backed performance with **99.9% agent availability** and **≤2s P95 latency** per agent call.
- **Autoscaling Infrastructure:** Dynamic GPU allocation per workload, with warm-started LLM pools to reduce cold-start delays.
- **Efficient Observability:** Token- and GPU-hour-level cost attribution without adding latency overhead.
- **Resiliency Features:** Circuit breakers, retry logic, fallback agents, and safe rollback mechanisms ensure graceful recovery.
- **Scalable Orchestration:** Multi-agent flows optimized with event-driven triggers and memory compaction for long-running tasks.



- **Drift & Anomaly Handling:** Proactive alerts for behavioral drift, hallucinations, or malicious tool calls, reducing error propagation in production.

## Tools and Azure Services Used

NexaStack integrates deeply with **Azure-native services**, ensuring seamless deployment for customers in Microsoft environments:

- **Compute & AI**
  - **Azure Kubernetes Service (AKS)** for containerized agent orchestration.
  - **Azure Machine Learning** for LLM hosting, model lifecycle, and versioning.
  - **Azure OpenAI Service** for GPT-based inference where applicable.
  - **Azure GPU VM SKUs (A100, H100, L40S, etc.)** for inference workloads.
- **Observability & Monitoring**
  - **Azure Monitor & Application Insights** for telemetry integration.
  - **LangFuse + OpenTelemetry** for AI-specific cognitive tracing.
- **Data & Storage**
  - **Azure Blob Storage / Data Lake** for logs, memory persistence, and audit trails.
  - **Azure Cosmos DB / Redis Cache** for real-time memory operations.
  - **Azure SQL Database** for agent metadata and orchestration state.
- **Security & Compliance**
  - **Azure Active Directory (Entra ID)** for IAM and RBAC enforcement.
  - **Azure Policy & Defender for Cloud** for compliance enforcement and threat detection.
  - **Key Vault** for secure secrets and credentials management.
- **Integration Services**
  - **Azure Logic Apps / Event Grid** for event-driven agent triggers.
  - **Service Bus** for inter-agent communication and workflow orchestration.
  - **Azure DevOps / GitHub Actions** for CI/CD pipelines of agent updates.

## Users of Agent

NexaStack agents serve multiple user personas inside enterprises:

- **Business Teams**
  - Finance teams for **cloud cost governance (FinOps)**.
  - Compliance teams for **regulatory monitoring and audit readiness**.
  - HR & Operations teams for **policy enforcement and automation**.
- **Technical Teams**
  - **SREs & DevOps Engineers** for incident RCA, drift detection, and CI/CD optimization.
  - **AI Governance Analysts** for bias monitoring and safe deployment checks.
  - **Developers** for testing new prompts, workflows, and orchestration patterns.
- **Executive & Leadership**
  - CIOs/CTOs for **cost transparency dashboards and SLA reports**.
  - Business leaders for **real-time insights from AI-driven analytics**.
- **Customer-Facing Roles**
  - Support teams using **AI agents for SLA-driven incident response**.
  - Sales/Service roles supported by **multi-agent orchestration** for queries and transactions.

## Dependencies

The successful deployment and operation of NexaStack depends on:

- **Cloud / Infra Dependencies**
  - Azure subscription with GPU-enabled VMs or AKS clusters.
  - Access to **Azure Monitor, Storage, and Security services**.
  - Network configuration (VNET, VPN, or ExpressRoute for on-premises hybrid).
- **Software Dependencies**
  - LangGraph / A2A for multi-agent orchestration.
  - LangFuse + OpenTelemetry for observability.
  - Redis / Vector DB (Pinecone, Milvus, or Azure Cognitive Search) for memory.
- **Data & Enterprise System Dependencies**
  - ERP/CRM systems (SAP, Oracle, Salesforce) for workflow integrations.
  - Data platforms (Snowflake, Databricks, Azure Synapse).
  - Identity providers (Azure AD, Okta) for role-based access control.
- **Organizational Dependencies**
  - Defined **SLA ownership model** for critical AI workloads.

- Adoption of **AI Governance policies** to align with compliance frameworks (GDPR, HIPAA, RBI, ISO 27001).
- Stakeholder buy-in from **Ops, Security, and Compliance teams** for rollout success.

## Key Benefits and Differentiators

NexaStack stands out as an **enterprise-ready AI MSP** by combining **agent lifecycle management, reliability, observability, and compliance** into a single unified platform.

### Key Benefits:

- **Reliable Operations** – SLA-backed uptime, auto-recovery, and safe rollback mechanisms.
- **Deep Observability** – Cognitive traces, token/GPU cost dashboards, and memory insights.
- **Secure Execution** – Sandbox environments, ACL-based controls, and prompt injection defense.
- **Compliance-First** – Continuous audit trails, lineage tracking, and bias monitoring for Responsible AI (RAI).
- **Scalability** – Multi-agent orchestration with support for complex workflows across industries.
- **Cost Transparency** – Token-level cost attribution and GPU-hour tracking for predictable TCO.

### Differentiators vs Traditional AI Platforms:

- Purpose-built **AgentOps & Infra Reliability** framework, not just model hosting.
- **Drift & anomaly detection** that proactively prevents unsafe outputs.
- **Multi-agent orchestration (A2A + LangGraph)** with traceability and governance.
- **MSP delivery model** with playbooks, RCA templates, and SLA-backed support.
- **Hybrid Deployment Flexibility** – deploy on Azure, multi-cloud, or on-premises GPU clusters.

# Value Proposition

NexaStack delivers **trustworthy, cost-efficient, and enterprise-grade AI agent operations** by acting as the **AI Managed Service Provider (MSP)** for organizations.

- For **CIOs & CTOs**: Assurance of **reliable AI adoption** with transparent cost models and enterprise SLAs.
- For **Ops & Security Teams**: Built-in **observability, rollback, and compliance guardrails** that reduce risk.
- For **Business Units**: **Accelerated time-to-value** with AI agents handling FinOps, compliance, DevOps, and customer operations.
- For **Developers & Data Teams**: A **sandboxed, orchestrated, and monitored environment** to safely innovate with AI agents.

In essence, NexaStack **de-risks AI adoption** by ensuring agents are **accountable, observable, secure, and auditable** — while scaling to meet enterprise needs.

## Conclusion

As enterprises transition from AI pilots to **production-scale agentic systems**, the lack of structured operational models creates risks in reliability, compliance, and cost management. NexaStack solves this gap by acting as an **AI Managed Service Provider**, embedding **AgentOps, observability, orchestration, anomaly detection, and governance** into enterprise AI ecosystems.

By partnering with NexaStack, organizations can:

- Confidently deploy multi-agent workflows in mission-critical environments.
- Achieve **audit-ready AI operations** aligned with compliance mandates.
- Ensure **high reliability and security** with continuous monitoring and rollback safety nets.
- Gain **transparent cost insights** for sustainable AI scaling.

**NexaStack is not just an AI platform — it is the operational backbone that ensures AI agents run securely, reliably, and responsibly at enterprise scale.**