# InferenceOps and management

Getting your first LLM into production is a big milestone. But staying there and scaling up requires more than just a working model. Without a reliable and standardized operational workflow, AI teams quickly find themselves in a maze of manual steps, patchwork tooling, and inconsistent processes.

This is where InferenceOps comes in: the practices and workflows that support ongoing model deployment, updates, and management at scale.

## Standardized deployment workflows

Every production LLM application should begin with clear, repeatable deployment processes. This helps avoid fire drills later on and ensures any engineer on the team can safely ship changes.

- **CI/CD pipelines for models**

  Just like traditional applications, models should be packaged, tested, and deployed automatically. A proper CI/CD setup ensures:

  - Changes are validated through test cases (e.g., regression, latency, token generation checks)

  - Infrastructure changes (e.g., resource requirements or caching configs) are reviewed alongside code

  - Deployment is repeatable and auditable

- **Release strategies: canary and blue-green deployments**

  Push models incrementally:

  - **Canary**: Route a small percentage of traffic to a new model version to monitor behavior before shifting all the traffic to it.

  - **Blue-green**: Keep two environments live (old and new) and switch traffic once the new version is verified. This minimizes downtime and rollback risk.

## Safe updates and fault tolerance

Once models are live, change becomes a constant. Whether it's performance tuning, bug fixes, or model swaps, your infrastructure must support safe iteration.

- **Rolling updates**. Deploy updates gradually across instances to avoid downtime. Each replica is replaced and verified before the next is rolled out.

- **Automatic rollback and alerting**. In case of failure (e.g., spike in latency, degraded accuracy, or traffic timeouts), the system should:

  - Alert engineers through monitoring dashboards or incident systems

  - Automatically revert to the previous model or routing configuration

  - Log the event for future auditing

- **Fault isolation**. Model failures should not bring down entire applications. Use retries, timeouts, circuit breakers, and load shedding to contain issues before they cascade.

## Centralized management at scale

What works for a few models quickly falls apart when you have dozens or hundreds, especially across different teams, cloud environments, and use cases.

- **Model registry and lifecycle tracking**. Ideally, you should maintain a central view of:
  - What models are deployed, where, and by whom
  - Version history and performance metrics
  - Ownership and compliance metadata
- **Unified control plane**. Deploy, monitor, and scale models from a single system across all clouds and environments. This eliminates siloed setups and reduces cross-team confusion.
- **Multi-region and multi-cloud support**. As you grow, your inference workloads may span multiple regions for latency, compliance, or failover. A unified deployment framework helps coordinate these rollouts and avoid drift between environments.

## Cost control and resource hygiene

Without proper InferenceOps, costs spiral and visibility disappears.

- **Idle GPU cleanup**. It's not uncommon for orphaned GPU instances to run for weeks, even months. Automate cleanup of unused or underutilized resources.
- **Access control and audit logs**. Ensure only authorized changes are made to production models, and that every deployment is logged for traceability.