## Choose Infrastructure (Cloud or On-Prem)

Begin by evaluating the operational paradigm most suitable for deploying your AI workloads. Public cloud environments—offered by vendors such as AWS, Azure, or Google Cloud—provide access to elastic compute, managed Kubernetes orchestration, and integrated storage and monitoring stacks. These platforms are highly adaptable and particularly advantageous in use cases that demand scalability, rapid provisioning, and minimal infrastructure maintenance. Alternatively, if data localization laws, compliance requirements, or latency-sensitive inference demand localized control, an on-premise deployment may be necessary. In this scenario, you'll assume complete responsibility for hardware, networking, and cluster management. If cloud deployment is selected, continue with Provision Cloud Infrastructure. If an on-premise approach is chosen, proceed to Set Up Bare-Metal or VMs.

## Provision Cloud Infrastructure

Within the cloud model, begin by provisioning GPU-accelerated instances capable of supporting high-throughput inference workloads. Deploy a container orchestration layer using managed Kubernetes services like EKS, AKS, or GKE. IAM roles must be tightly scoped, virtual networks properly segmented, and storage classes defined to support both ephemeral and persistent data flows. Leverage Terraform or Pulumi to codify your infrastructure so that environments remain reproducible and auditable. Configure telemetry agents, centralized logging, and secrets management integrations (e.g., AWS Secrets Manager or Azure Key Vault) to ensure compliance with security and observability baselines. Once your cloud control plane and supporting services are in place, continue with Select Model Serving Platform.

## Set Up Bare-Metal or VMs

For on-premise deployment, provision physical servers or virtual machines equipped with compatible GPUs such as the NVIDIA A100 or H100. Container runtime environments like Docker or CRI-O must be installed and configured, and Kubernetes orchestration should be bootstrapped via kubeadm, Rancher, or equivalent tooling. Networking must be explicitly defined using static IP assignment, internal DNS, and possibly VLAN tagging. High-throughput local storage—typically NVMe-based—should be mounted using LVM or RAID strategies. Security measures such as firewall zoning, VPN tunneling, and device-level hardening are essential. Following this foundational setup, continue with Select Model Serving Platform.

## Select Model Serving Platform

Next, identify a model inference engine compatible with your deployment goals and supported model formats. NVIDIA Triton Inference Server, TensorFlow Serving, or TorchServe are among the preferred options, especially where dynamic batching, GPU multiplexing, or multi-model deployment is required. Encapsulate the model, its preprocessing pipeline, and health-check interfaces within a container. Define exposure methods via REST or gRPC depending on downstream system compatibility. Helm or Kustomize can be employed to template deployments and ensure repeatability across environments. After your model serving platform is prepared, proceed to Deploy Inference Controller.

## Deploy Inference Controller

Deploy the model-serving containers into your orchestration fabric using StatefulSets or Deployments, depending on whether session affinity or persistence is needed. Set up

ConfigMaps for runtime parameters and Kubernetes Secrets for sensitive credentials. Attach PersistentVolumeClaims if model files are dynamically pulled or logs are retained. Use Horizontal Pod Autoscalers (HPA) or KEDA to scale replicas based on system load or request frequency. Apply node selectors, tolerations, and affinity rules to bind workloads to GPU-enabled nodes. Include readiness and liveness probes to support high availability. Once the inference infrastructure is deployed, branch immediately into parallel configuration of the following components: Set Up Input Channels, Set Up Output & Storage Systems, Install Metrics & Addons, and Configure RBAC & Security.

## Set Up Input Channels

Enable ingress paths that feed real-time or batch inputs into your inference service. This may involve HTTP endpoints, WebSocket streams, or message queue subscribers using Kafka or RabbitMQ. Input schemas must be enforced using validation libraries, and authentication mechanisms such as OAuth2 or JWT must be implemented at entry points. Leverage API gateways or mesh ingress controllers like Istio Gateway or Envoy Proxy to manage routing, retries, and rate limiting. Ensure that trace and correlation IDs are attached to each inbound request for auditability and observability. Because these entry points are externally exposed, ensure that their configuration is tightly coupled with the RBAC and security policies defined in Configure RBAC & Security to maintain a secure interface surface.

## Set Up Output & Storage Systems

Configure output pipelines to persist prediction results or forward them downstream. Structured outputs may be directed to relational databases like PostgreSQL, while unstructured blobs or batch exports can be sent to object stores like Amazon S3 or Azure Blob. Redis or Memcached can support short-lived caching of responses. Ensure encryption at rest for all output data and enable TLS across all outbound interfaces. Streaming outputs to data pipelines or dashboards may require CDC mechanisms or dedicated connectors. Integrate output metadata with data catalogs for lineage tracking and auditing. Once output storage is operational, the system is ready to accept and persist results from live inference traffic.

## Install Metrics & Addons

Deploy observability tooling that enables visibility into all operational aspects of the system. Prometheus should be configured to scrape performance and system metrics, while Grafana visualizes latency and throughput. Fluent Bit or Fluentd can route logs to Elasticsearch, Loki, or a central SIEM platform. Use OpenTelemetry to instrument distributed traces and correlate service spans. Define dashboards that reflect system health, SLOs, and capacity planning trends. Set up alerts for anomalies, errors, and degradation, notifying responders via channels such as Slack, PagerDuty, or Opsgenie. Once observability components are in place and verified, continue to Validate the AI System to monitor system behavior under realistic conditions.

## Configure RBAC & Security

Apply RBAC policies within Kubernetes to restrict access to APIs, secrets, and runtime logs. Enforce mTLS between services to prevent eavesdropping and service spoofing. Use ingress controllers to filter inbound IPs and limit request rates. Secrets must be managed with tools like HashiCorp Vault, Sealed Secrets, or external secrets operators. Perform static and dynamic scans of containers, and enforce AppArmor or SELinux profiles on nodes to restrict privileged operations. Where applicable, implement DLP scanners to

prevent exfiltration of sensitive data through inputs or outputs. Once these security controls are fully applied, proceed to Validate the AI System to finalize the secure deployment perimeter.

## Validate the AI System
Final system validation includes stress testing, behavioral verification, and resiliency analysis. Tools like Locust, Artillery, or k6 simulate user demand and reveal performance boundaries. Confirm that autoscaling logic triggers correctly and that observability data reflects real-time system conditions. Authentication and authorization flows must be re-tested under load. End-to-end trace analysis should be conducted to ensure no bottlenecks exist across the inference path. Outputs are validated against test datasets or live sampling tools. If validation exposes critical issues—such as misconfigurations, failed routing, or unmet performance thresholds—you must rollback to Deploy Inference Controller, revise deployment manifests, and re-execute the deployment iteration. Upon successful validation, the system may be marked production-ready and released under a monitored rollout policy.