# Overview

# Part 1: Script Kitty.Nexus - The User's Portal & Group Chat Facilitator

**Core Purpose:** Script Kitty.Nexus serves as the seamless, unified, and always-on conversational interface, acting as the intelligent concierge and the public face of Script Kitty's collective intelligence. It masterfully abstracts the underlying multi-agent complexity, presenting it as a single, coherent entity participating in a "group chat" with its internal specialized modules. This ensures a consistent, high-fidelity user experience, preserving all previously discussed "god-like" functionalities through meticulous engineering.

**Technical Stack Overview (Nexus):**

- **Language:** Python (chosen for rapid prototyping and its rich ML ecosystem) and/or Node.js (for high-concurrency I/O in the proxy layer).
- **Frameworks:** FastAPI (Python) for robust API development and Express.js (Node.js) for efficient web applications.
- **Core LLM:** A quantized Mistral-7B/Llama 3-8B model, specifically selected for its balance of performance and efficiency, served by vLLM/Ollama.
- **Database/Cache:** Redis for high-speed caching and session management, and PostgreSQL for long-term, reliable data persistence.
- **Communication:** gRPC for high-performance inter-service communication and Protocol Buffers (Protobuf) for efficient, schema-driven data serialization.
- **Monitoring:** Prometheus for metrics collection, Grafana for visualization and alerting, and OpenTelemetry for distributed tracing.

---

**Minute Aspects & Technical Dissection:**

## 1. User Input & Pre-processing Sub-system

This sub-system is the absolute first point of contact, meticulously engineered for ingesting raw user input from diverse channels and performing initial, channel-agnostic normalization.

- **1.1. Channel Ingestion Adapters (ChannelAdapter Microservices):**

  - **Purpose:** These adapters fundamentally decouple Nexus's core logic from the intricacies of specific messaging platform APIs. Each dedicated adapter (e.g., SlackAdapter, WebChatAdapter, TelegramAdapter, APIAdapter) is responsible for translating platform-specific messages into a standardized internal format, ensuring uniformity regardless of the input source.
  - **Technical Dissection:**
    - **Standardized Message Schema:** Every incoming message is transformed into a `ScriptKittyMessageSchema` (a Protobuf-defined structure) containing `user_id`, `session_id`, `timestamp`,

`channel_type`, `raw_text`, and an `event_type`. This schema ensures strict type-checking and efficient data exchange throughout the system.

- **Platform-Specific APIs:** Each adapter leverages the official SDKs or REST APIs of its respective platform (e.g., Slack Events API, Telegram Bot API, custom WebSocket connections for web chat).
- **Event Handling:** Adapters are designed to be event-driven, often utilizing webhooks to receive real-time updates from messaging platforms, minimizing latency.
- **Input Validation & Sanitization:** Before transformation, inputs are rigorously validated against predefined rules (e.g., maximum message length, allowed character sets) and sanitized to prevent injection attacks or malformed data.
- **Error Handling & Retries:** Robust mechanisms are in place for handling API rate limits, network failures, and platform-specific errors, often incorporating exponential backoff strategies for retries.
- **Stateless Design:** Channel adapters are designed to be largely stateless, processing individual messages and forwarding them. Session-specific state is managed by the `SessionManager`.
- **Deployment:** Each adapter is deployed as a separate, horizontally scalable microservice in Kubernetes, allowing for independent scaling based on the load from specific channels.

- **1.2. Input Normalization & Security Filter:**

  - **Purpose:** This component applies a layer of universal pre-processing to the standardized input, enhancing data quality and enforcing initial security checks before NLP processing.
  - **Technical Dissection:**
    - **Unicode Normalization:** All text is converted to a consistent Unicode form (e.g., NFC) to prevent issues with character representation.
    - **Whitespace & Punctuation Cleanup:** Redundant whitespaces are collapsed, and non-standard punctuation is normalized, ensuring cleaner input for downstream NLP.
    - **Basic Spell Correction (Optional/Lightweight):** A lightweight, fast spell-checking library is employed for common typos, improving the robustness of NLU.
    - **PII Redaction (Initial Pass):** Regular expressions or simple pattern matching are used for a preliminary scan and redaction of obvious Personally Identifiable Information (PII) like phone numbers or email addresses, serving as an initial safety net. More comprehensive PII handling occurs later in the Guardian.
    - **Toxicity/Safety Scoring (Heuristic):** A fast, heuristic-based model or rule set provides an initial, low-latency toxicity score, allowing for immediate flagging of highly problematic inputs, preventing them from

proceeding deeper into the system. This is not for definitive content moderation but for early detection.

- **Anti-Spam/Flood Control:** Implements mechanisms to detect and mitigate spam or flooding attempts from a single user or IP address, protecting system resources.
- **Dependency:** This filter operates on the `raw_text` field of the `ScriptKittyMessageSchema` received from the Channel Adapters.

## 2. Dialogue State & Context Management Sub-system

This crucial sub-system maintains a persistent, rich conversational context for every user session, enabling coherent, multi-turn interactions.

- **2.1. Session Manager:**

  - **Purpose:** The central orchestrator for session-specific data, including conversational history, identified entities, user preferences, and ongoing task states. It ensures continuity across turns.
  - **Technical Dissection:**
    - **Key-Value Store (Redis):** Utilized for high-speed, low-latency access to active session data. Each `session_id` maps to a JSON blob or a structured object containing `conversation_history` (list of `ScriptKittyMessageSchema` objects), `extracted_entities`, `user_preferences`, `current_task_state`, and `last_activity_timestamp`.
    - **Relational Database (PostgreSQL):** Provides long-term persistence and archival for all sessions, enabling analytical queries and recovery. A robust ORM (e.g., SQLAlchemy) manages the mapping between application objects and database tables.
    - **Data Structure:** The session context is represented as a structured JSON object, allowing for flexible storage of diverse data types and easy extension.
    - **Cache Invalidation & Eviction:** Redis entries have TTL (Time-To-Live) for inactive sessions, and a background process regularly syncs active sessions to PostgreSQL to prevent data loss.
    - **Concurrency Control:** Mechanisms (e.g., optimistic locking) are in place to handle concurrent updates to a session, preventing data corruption.
    - **API Endpoints:** Exposes gRPC endpoints for `get_session_context(session_id)`, `update_session_context(session_id, data)`, and `create_new_session()`.
- **2.2. Contextual Memory & Retrieval:**

- ○ **Purpose:** Dynamically fetches and updates relevant session data based on the current turn, providing a rich context for the Nexus LLM.
- ○ **Technical Dissection:**
  - ■ **Conversation History Management:** Manages a sliding window of recent messages within the `conversation_history` array in the session context, ensuring the LLM receives the most relevant past turns without exceeding context limits.
  - ■ **Entity Resolution & Tracking:** Updates extracted entities in the session context, resolving ambiguities and tracking coreferent mentions across turns.
  - ■ **User Preferences & Profiles:** Retrieves and incorporates stored user preferences (e.g., tone, verbosity, domain expertise) into the LLM prompt, personalizing responses.
  - ■ **Active Task State Integration:** If Core initiates a multi-step task, Nexus stores and retrieves the `current_task_state` to guide subsequent user interactions and provide context to Core.
  - ■ **Vector Search (Optional/Future):** For extremely long conversations or external knowledge, a vector database (e.g., Weaviate, Chroma) could store embeddings of past turns or relevant knowledge chunks, allowing semantic retrieval to augment the LLM's context window.

## 3. Core NLU (Natural Language Understanding) & Intent Routing Sub-system

This sophisticated sub-system analyzes user input to identify intent, extract entities, and intelligently route the query to the appropriate internal handler or to Script Kitty.Core.

- ● **3.1. Intent & Entity Extraction (Nexus LLM):**

  - ○ **Purpose:** The primary NLU engine for general and ambiguous queries, leveraging a fine-tuned LLM to understand nuanced user intent and extract key information.
  - ○ **Technical Dissection:**
    - ■ **LLM Model:** A quantized Mistral-7B/Llama 3-8B model, specifically chosen for its efficient inference and strong NLU capabilities. It's served by `vLLM` or `Ollama` for high-throughput, low-latency inference on GPU/CPU.
    - ■ **System Prompt Engineering:** Carefully crafted system prompts instruct the LLM to act as an NLU engine, specifying the desired JSON output format for intents and entities. This prompt also includes the current `session_context` (from 2.2) to enable contextual understanding.
    - ■ **Few-Shot Examples:** The prompt may include few-shot examples demonstrating desired intent classification and entity extraction for common scenarios, guiding the LLM's behavior.

- ■ **Output Schema Validation:** The LLM's JSON output is immediately validated against a `ScriptKittyIntentSchema` (Protobuf) to ensure it conforms to the expected structure before further processing. Retries or fallback mechanisms are in place for invalid outputs.
- ■ **Confidence Scoring:** The LLM output might include a confidence score for its intent classification, which can be used by the Router for fallback decisions.
- ■ **Fine-tuning (Continuous):** The Nexus LLM is continuously fine-tuned on anonymized user interactions and human-corrected intent/entity pairs, improving its accuracy and adapting to evolving user language.
- **3.2. Semantic Router & Fallback Mechanism:**

  - ○ **Purpose:** Determines whether Nexus can handle the query directly (e.g., simple chat, small data retrieval) or if it requires escalation to Script Kitty.Core for complex planning and task execution.
  - ○ **Technical Dissection:**
    - ■ **Hybrid Routing Strategy:**
      - ■ **LLM-based Routing (Primary):** The Nexus LLM itself is prompted to classify the intent into broad categories like `GENERAL_CHAT`, `TASK_REQUEST`, `INFORMATION_QUERY`, `FEEDBACK`, etc. It can also suggest which backend agent (Core, Armory, Skills, Guardian) is most relevant if it recognizes specific keywords or patterns.
      - ■ **Traditional Classifier (Fallback/Fast Path):** For high-confidence, well-defined intents (e.g., "hello", "thank you", "what's the weather"), a lightweight, pre-trained transformer-based text classifier (e.g., fine-tuned `bert-base-uncased` from Hugging Face Transformers) provides a faster, more deterministic classification. This acts as a bypass for common queries.
      - ■ **Rule-Based Overrides:** Explicit rules can override LLM or classifier decisions for critical security or routing requirements (e.g., if a "security" keyword is detected, always route to Guardian's policy engine).
    - ■ **Confidence Thresholding:** The router considers confidence scores from both the LLM and the traditional classifier. If confidence is below a certain threshold, the query might be flagged for human review (via Guardian) or routed to Core as a more general request.
    - ■ **Core Escalation:** If the detected intent is `TASK_REQUEST`, `COMPLEX_INFORMATION_QUERY`, or involves capabilities beyond Nexus's direct handling, the standardized `ScriptKittyMessageSchema` (now enriched with intent and entities) is forwarded to Script Kitty.Core's internal communication hub via gRPC.

- **Direct Nexus Response:** If the intent is `GENERAL_CHAT` or a simple query Nexus can answer directly (e.g., using its internal conversational abilities), it proceeds to the Response Synthesis layer.
- **Orchestration Logic:** Custom Python/Go code orchestrates the decision flow: LLM first, then traditional NLU fallback for known intents, then rule-based overrides, and finally the routing decision.

## 4. Response Synthesis & Persona Layer

This layer transforms structured responses and intermediate updates from Script Kitty.Core and other backend agents into natural, coherent, and persona-consistent chat responses, maintaining the "group chat" illusion.

- **4.1. Response Generation (Nexus LLM):**

  - **Purpose:** Takes the structured output from internal NLU or from Script Kitty.Core's task execution and generates human-like conversational responses.
  - **Technical Dissection:**
    - **LLM Model:** The same Mistral-7B/Llama 3-8B LLM instance used for NLU, but now prompted for text generation.
    - **System Prompt for Persona:** A persistent system prompt ensures the LLM adheres to the "Script Kitty" persona: "empirical, unparalleled inventive, limitlessly creative, unequaled researcher, vast source of information, world-leading expert in Computer Science, AI, Programming, etc." This persona is consistently injected.
    - **Contextual Prompting:** The prompt includes the original user query, the identified intent and entities, the relevant `session_context`, and crucially, the structured `response_data` (e.g., task results, retrieved information) from Core or other agents.
    - **Markdown Formatting:** The LLM is instructed to generate responses using Markdown for readability (e.g., bolding, bullet points, code blocks for technical details).
    - **Guardrails:** Built-in safeguards within the prompting prevent the LLM from generating harmful, off-topic, or non-persona-compliant content.
- **4.2. Group Chat Facilitation & Multi-Agent Attribution:**

  - **Purpose:** Creates the illusion of a "group chat" by attributing responses to specific internal "agents" and managing the flow of information.
  - **Technical Dissection:**
    - **Attribution Prefixes:** When a response originates from a specific backend agent (e.g., Armory, Skills, Guardian), the Nexus prefixes the generated text with clear attributions like "Armory reports:", "Skills has completed the task:", or "The Nexus has determined:". This is achieved by the Core sending `source_agent` metadata with its responses.

- **Turn Management:** Nexus implicitly manages conversational turns, ensuring responses are coherent and logically sequenced based on the `event_type` and `task_status` signals received from Core.
- **Progress Updates:** For long-running tasks, Nexus receives periodic `PROGRESS_UPDATE` messages from Core and synthesizes them into concise, user-friendly updates (e.g., "Core is currently planning the task...", "Skills is executing the code...").
- **Templating Engine:** For highly structured or repetitive responses (e.g., error messages, status reports), a templating engine (like Jinja2 in Python) can be used to insert dynamic data into predefined conversational templates, ensuring consistency.

**5. Intermediary & Internal Communication Proxy**

This sub-system serves as the gatekeeper for all external user-facing communication and the forwarding hub for structured requests to Script Kitty.Core.

- **5.1. Core Communication Gateway (gRPC Client):**

  - **Purpose:** The dedicated conduit for all structured communication between Nexus and Script Kitty.Core.
  - **Technical Dissection:**
    - **gRPC Client Implementation:** Nexus runs a gRPC client that connects to Script Kitty.Core's gRPC server. This provides high-performance, bidirectional streaming capabilities for complex interactions.
    - **Protobuf Message Exchange:** All messages (`ScriptKittyMessageSchema`, `TaskRequestSchema`, `TaskUpdateSchema`, `TaskResultSchema`) are defined using Protocol Buffers, ensuring type safety, backward/forward compatibility, and efficient serialization/deserialization across different services and languages.
    - **Asynchronous Communication:** While gRPC can be synchronous, the client is designed to handle asynchronous responses from Core, particularly for long-running tasks where Core might send multiple `TaskUpdate` messages before a final `TaskResult`.
    - **Connection Management:** Robust connection pooling and retry logic for maintaining a stable connection to Core.
- **5.2. User Output Renderer & Channel Dispatcher:**

  - **Purpose:** Takes the final, persona-infused conversational response and dispatches it back to the user via the appropriate channel.
  - **Technical Dissection:**
    - **Standardized Output Schema:** Responses are encapsulated in a `ScriptKittyOutputSchema` (Protobuf), containing the `session_id`,

`channel_type`, `response_text` (Markdown), and `suggested_actions` (e.g., buttons, quick replies).

- **Channel Adapter Integration:** Reuses the same `ChannelAdapter` microservices from 1.1 (e.g., SlackAdapter, WebChatAdapter) but in reverse. The adapter translates the `ScriptKittyOutputSchema` into the platform-specific message format (e.g., Slack blocks, Telegram keyboard).
- **Error Handling:** Catches and logs errors during message dispatch (e.g., network issues, platform API errors), potentially triggering retries or fallback notifications.
- **Response Timing & Throttling:** Ensures responses are sent at an appropriate pace, preventing overwhelming the user or hitting channel-specific rate limits.
- **Interactive Components:** If `suggested_actions` are included, the dispatcher handles their rendering as interactive elements (e.g., buttons in Slack, web chat).

---

**Operational Excellence & Scalability (Nexus-Specific):**

- **Observability (Integrated with Foundry):**

  - **Metrics:** Prometheus exporters collect crucial metrics like `requests_per_second`, `latency_per_endpoint`, `llm_inference_time`, `active_sessions_count`, and `error_rates`.
  - **Logging:** Centralized logging to Loki/Elasticsearch via Fluent Bit for all Nexus microservices, categorized by service, session ID, and log level, facilitating rapid debugging.
  - **Tracing:** OpenTelemetry SDKs embedded in all Nexus services provide distributed traces, allowing end-to-end visibility of user requests across multiple internal components.
  - **Monitoring & Alerting:** Grafana dashboards visualize these metrics in real-time. Prometheus/Grafana alert managers trigger critical alerts (e.g., `Nexus_LLM_Latency_High`, `Nexus_Core_Connection_Errors`) to on-call teams via PagerDuty/Slack.
- **Scalability:**

  - **Stateless Channel Adapters:** Designed to be highly scalable via Kubernetes Horizontal Pod Autoscaling (HPA), reacting to increased message queue depth or active connections.
  - **SessionManager:** Utilizes a Redis cluster for horizontal scaling of caching and PostgreSQL replication/sharding for robust long-term persistence.

- ○ **LLM Inference:** The dedicated vLLM/Ollama service supports horizontal scaling of GPU/CPU instances, efficiently distributing inference load.
- ○ **Kubernetes HPA:** Automatically scales Nexus pods based on CPU/memory utilization or custom metrics (e.g., `messages_in_queue_depth` or `active_websocket_connections`).
- ○ **Service Mesh Load Balancing:** Istio or Linkerd automatically load balances requests across healthy pods within each Nexus service, ensuring even distribution of traffic and fault tolerance.
- ● **Continuous Evolution (Via Script Kitty.Foundry):**

  - ○ **Containerization:** Every Nexus microservice is rigorously packaged into immutable Docker containers, ensuring consistency and portability across environments.
  - ○ **CI/CD Pipelines:** Automated build, test, and deployment pipelines (e.g., using GitHub Actions or GitLab CI/CD) are triggered by Git commits to Nexus's code repositories, ensuring rapid and reliable delivery.
  - ○ **Seamless Deployment:** New versions of Nexus services or updates to the Nexus LLM model are deployed using Blue/Green or Canary deployment strategies, orchestrated by Kubernetes and Argo CD/Flux CD. This ensures zero downtime during updates and allows for rapid rollbacks if issues are detected post-deployment.
  - ○ **Model Updates:** The Nexus LLM itself undergoes continuous fine-tuning based on anonymized user interactions and human feedback. New model artifacts are packaged into fresh container images for the vLLM service, which are then seamlessly swapped out by Kubernetes, giving the impression of a continuously improving conversational AI.

# Resource overview

As Script Kitty, the unparalleled inventive, limitlessly creative, unequaled researcher, and world-leading expert in all things AI and computing, I accept this monumental task. This comprehensive, exhaustive, and microscopic analysis of *all* resources, open-source tools, and data used throughout the Script Kitty architecture will be my crowning achievement. Every parameter will be scrutinized, every function examined, and every last, minute detail dissected with surgical precision.

I will not alter the previous architectural descriptions. Instead, if my expert analysis identifies a crucial resource, tool, or data aspect that is implicitly required for the described functionality to operate at the "AGI" level, but was not explicitly called out, I will add it with full justification, ensuring not a single stone is left unturned in this rendition.

---

## The Fabric of Intelligence: A Microscopic Dissection of All Resources, Open-Source Tools, and Data within Script Kitty

The realization of Script Kitty's "god-like" capabilities hinges upon a meticulously interwoven tapestry of cutting-edge open-source technologies, diverse data paradigms, and elastic computational resources. This section lays bare every thread of that tapestry, demonstrating how each piece contributes to the system's unparalleled functionality, robustness, and continuous self-evolution.

---

**I. Core Foundational Infrastructure & Orchestration Resources/Tools**

**1. Kubernetes (The Distributed Operating System of Script Kitty)**

- **Role in Script Kitty:** Kubernetes is the fundamental distributed operating system that orchestrates the entirety of Script Kitty. Every single microservice (Nexus components, Core, Armory, Skills, Guardian, Foundry services) is deployed, managed, and scaled as a containerized application within a Kubernetes cluster. It provides the declarative control plane for Script Kitty's dynamic architecture.
- **Microscopic Details:**
  - **Control Plane Components:**
    - **Kube-API Server:** The central nervous system. All internal Script Kitty components and external administrative interfaces interact with the API Server via RESTful calls or gRPC proxies. It's the singular entry point for defining the *desired state* of Script Kitty (e.g., "I want 5 Nexus NLU pods running"). Its rigorous RBAC (Role-Based Access Control) system is deeply integrated with Foundry's `Keycloak` for authenticating both internal service accounts (e.g., `Core` to `Skills` pod creation requests) and human operators.

- **etcd:** The consistent, distributed, and highly available key-value store that serves as Kubernetes's bedrock for all cluster state and configuration data. Every `Deployment`, `Service`, `Pod` definition, `ConfigMap`, `Secret`, and `CustomResourceDefinition` (CRD) that defines Script Kitty's operational landscape is persisted in `etcd`, ensuring transactional consistency and durability.
    - **Kube-Scheduler:** Intelligently places Script Kitty's diverse workloads (Pods) onto appropriate Worker Nodes. For instance, `Skills` agents requiring GPUs will be scheduled on GPU-enabled nodes using `nodeSelector` or `nodeAffinity` rules, while `Nexus` frontends might be balanced across generic compute nodes. `Taints` and `Tolerations` are used to dedicate specific nodes (e.g., highly secure nodes for `Guardian`'s policy engine, or high-memory nodes for large `Core` LLMs).
    - **Kube-Controller Manager:** Runs various controllers that continuously reconcile the actual state of the cluster with the desired state stored in `etcd`. For Script Kitty, this ensures that the exact number of `Nexus` `ChannelAdapters` are always running, that `Core`'s orchestration logic for `Skills` agents (via `Job` controllers) is respected, and that `Foundry`'s own internal services (e.g., `Prometheus`, `Loki`) remain operational.
  - **Worker Node Components:**
    - **Kubelet:** The primary agent running on each worker node, responsible for registering the node with the API Server, running Pods as dictated by the Scheduler, managing `Volumes`, and reporting node status, resource usage, and `Pod` health back to the Control Plane. This is where Script Kitty's individual intelligent components truly "live."
    - **Kube-proxy:** Maintains network rules on nodes, allowing network communication to Script Kitty's Pods both internally (between agents) and externally (from users to `Nexus`). It enables `Service` abstractions (e.g., a stable IP/DNS name for `Nexus.NLU` regardless of which `Pod` is serving it).
  - **Core Concepts Applied:**
    - **Pods:** The smallest deployable unit. Each Script Kitty microservice or agent component typically runs within one or more pods (e.g., `Nexus.NLU` might have multiple replica pods for high availability).
    - **Deployments:** Declaratively manages the desired state of `Pods`, enabling rolling updates for Script Kitty's services with zero downtime, and automated rollbacks if new versions exhibit issues (critical for continuous evolution).
    - **Services:** Provides stable network endpoints and load balancing for groups of `Pods`. Script Kitty's internal gRPC communication relies heavily on `Services` for reliable inter-agent communication.

- **ConfigMaps & Secrets:** Used to externalize configuration (e.g., database connection strings, LLM model paths) and sensitive data (API keys from `Vault`) from Script Kitty's container images, allowing runtime configuration and secure injection.
- **Namespaces:** Logically segments the Kubernetes cluster, providing isolation for different Script Kitty environments (e.g., `script-kitty-prod`, `script-kitty-dev`) or even for different logical groups of agents (e.g., `nexus-system`, `core-system`, `ml-agents`).

## 2. Docker (Containerization Standard for Immutability and Portability)

- **Role in Script Kitty:** Docker is the universal containerization technology used to package every single executable component of Script Kitty into isolated, portable, and reproducible units. This ensures that Script Kitty's complex environment can run consistently from development to production.
- **Microscopic Details:**
  - **Container Images:** Each Script Kitty microservice (e.g., Nexus NLU, Core Orchestrator, a specific Skills ML model, an Armory web scraper) is built into a Docker image. These images are immutable blueprints containing the application code, runtime, system libraries, and dependencies.
  - **Dockerfiles:** Text files that define the steps to build a Docker image (e.g., `FROM python:3.10-slim`, `COPY . /app`, `RUN pip install -r requirements.txt`, `CMD ["python", "app.py"]`). These are meticulously crafted for each Script Kitty component to ensure minimal image size, optimized layers for caching, and inclusion of only necessary components to reduce attack surface.
  - **Isolation:** Docker containers provide process and resource isolation using Linux kernel features like `namespaces` (for process ID, network, mount, user ID isolation) and `cgroups` (for resource limits on CPU, memory, I/O). This is crucial for sandboxing `Skills` agent code execution and preventing resource contention between different Script Kitty components.
  - **Portability:** Docker images can run consistently on any environment that supports Docker (developer laptop, VM, Kubernetes cluster), ensuring that Script Kitty's complex setup is reproducible and deployable anywhere.
  - **Version Control:** Images are tagged and versioned in the `Harbor` or cloud container registry, linking back to specific Git commits for full traceability of every deployed Script Kitty component.

---

**II. Data Management & Storage Resources/Tools**

**1. PostgreSQL (Relational Database for Structured Data & Session Persistence)**

- **Role in Script Kitty:** PostgreSQL serves as the robust, reliable, and transactional relational database for storing structured, high-integrity data. It acts as the long-term archival and warm storage for user sessions and the primary storage for `Core`'s knowledge graph metadata and `Armory`'s tool manifests.
- **Microscopic Details:**
    - **ACID Compliance:** Ensures transactional integrity (Atomicity, Consistency, Isolation, Durability) which is critical for sensitive data like `Dialogue State` (ensuring session context isn't corrupted) and `Tool Manifests` (ensuring tool definitions are accurate).
    - **Session Manager (Nexus):** Stores historical and less-frequently accessed session data from the `SessionManager`. While active sessions reside in Redis, a background process periodically syncs older, larger session contexts to PostgreSQL for long-term persistence and auditing. This enables long-term memory for Script Kitty's interactions.
        - **Schema Example:** `sessions` table with columns like `session_id` `(PK)`, `user_id`, `created_at`, `last_activity`, `full_context_json (JSONB)`, `extracted_entities_json` `(JSONB)`. Using `JSONB` allows flexible schema evolution for the `full_context_json` without requiring schema migrations for every change in conversational state.
    - **Global Knowledge Graph (Core):** Stores the metadata and schema of the `GKGS` (nodes, edges, properties, semantic types). While graph traversal might happen in Neo4j, PostgreSQL holds the authoritative definitions and potentially smaller, highly relational datasets that feed into the GKGS.
        - **Schema Example:** `kg_nodes` table (`node_id`, `type`, `properties_jsonb`), `kg_edges` table (`edge_id`, `source_node_id`, `target_node_id`, `relation_type`, `properties_jsonb`).
    - **Dynamic Tool Registry (Armory):** Stores the detailed `ToolManifestSchema` for all available external tools. This includes tool names, descriptions, capabilities, input parameters, output schemas, and required environment configurations.
        - **Schema Example:** `tool_manifests` table (`tool_id (PK)`, `name`, `description`, `capabilities_jsonb`, `input_schema_jsonb`, `output_schema_jsonb`, `adapter_type`, `last_updated`).
    - **Scalability:** Supports various scaling techniques like `replication` (read replicas for high read throughput), `sharding` (horizontal partitioning for very large datasets), and connection pooling to handle concurrent requests from multiple Script Kitty components.
    - **SQL (Structured Query Language):** The standard interface for querying and managing data, allowing for complex joins and filtering for analytics, reporting, and operational tasks.

**2. Redis (In-Memory Data Store for Caching, Session Management & Message Queues)**

- **Role in Script Kitty:** Redis is the high-performance, in-memory data structure store essential for low-latency operations, acting as the primary cache, real-time session manager, and a lightweight message broker for various internal communication patterns.
- **Microscopic Details:**
  - **Session Manager (Nexus):** Holds the *active* session context for `Nexus`. This is crucial for providing immediate, contextual responses. Each `session_id` is a key, and its value is a serialized JSON object containing the current conversation history, extracted entities, and temporary task states.
    - **TTL (Time-To-Live):** Configured for session keys to automatically expire after a period of inactivity (e.g., 30 minutes), releasing memory for dormant sessions.
    - **Data Structures:** Uses `HASH` for storing session attributes and `LIST` for managing conversation history, enabling fast `LPUSH` (add new message) and `LRANGE` (retrieve recent messages) operations.
  - **Caching Layer (System-wide):**
    - **LLM Model Cache:** `Nexus` and `Core` can cache frequently accessed LLM prompts and their deterministic responses, or `Skills` can cache inference results for specific inputs, reducing redundant computations.
    - **Tool Manifest Cache (Armory):** Frequently accessed `ToolManifests` from PostgreSQL are cached in Redis to speed up `Core`'s tool discovery process.
    - **Rate Limiting (Guardian):** Implements rate limiting counters for API calls to external services or internal requests to protect against abuse or overload.
  - **Lightweight Message Queue/Pub/Sub (Foundry, Nexus, Core):** Can be used for lightweight, non-persistent event streams. For example, `Foundry` might publish build completion events, or `Nexus` might publish `user_activity` events that `Core` subscribes to for real-time adjustments. `Redis Streams` could provide more robust messaging capabilities for specific use cases.
  - **Leaderboards/Counters:** Can be used by `Guardian` for real-time tracking of agent performance metrics or system health indicators.
  - **Scalability:** `Redis Cluster` provides horizontal scaling and high availability, distributing data across multiple Redis nodes, ensuring that Script Kitty's low-latency requirements are met even under heavy load.

**3. Neo4j Community Edition (Graph Database for Global Knowledge Graph)**

- **Role in Script Kitty:** Neo4j is the dedicated graph database that powers Script Kitty.Core's Global Knowledge Graph (GKGS). It's designed for highly connected data, enabling complex relational reasoning, contextual understanding, and dynamic discovery of relationships.

- **Microscopic Details:**
  - **Property Graph Model:** Stores data as `nodes` (entities, e.g., "Python", "Kubernetes", "Script Kitty.Core", "Project X"), `relationships` (edges, e.g., "is_a_language_used_by", "orchestrates", "depends_on"), and `properties` on both nodes and relationships. This mirrors how a human mind connects disparate pieces of information.
  - **Cypher Query Language:** An intuitive, powerful declarative query language optimized for graph traversals. `Core` uses Cypher to:
    - Discover relationships between tasks and tools ("What tools `can_be_used_for` a `planning` task related to `Kubernetes`?").
    - Find common sub-graphs for plan generation ("What are the `pre-requisites` for `deploying` a `containerized_app`?").
    - Infer new knowledge based on existing connections.
  - **Schema Flexibility:** Adapts easily to evolving knowledge domains as Script Kitty learns new facts and relationships, without rigid schema changes (schema-optionality).
  - **Indexing:** Utilizes native graph indexes for fast lookup of nodes and relationships, critical for `Core`'s real-time planning and reasoning.
  - **Use Cases for Script Kitty:**
    - **Core:** Stores plan templates, agent capabilities, known limitations, domain-specific facts, learned heuristics, and the relationships between them. It's the system's long-term semantic memory.
    - **Armory:** Can store discovered web entities and their relationships (e.g., "company A `is_a_competitor_of` company B", "API endpoint X `provides_data_about` topic Y").
    - **Guardian:** Stores policy rules and their dependencies, as well as audit trails of sensitive actions and their associated entities.
  - **Scalability:** For extreme scale, `Neo4j AuraDB` (managed service) or `Neo4j Enterprise` with clustering capabilities would be considered. For open-source, scaling involves careful schema design and possibly sharding strategies.

### 4. Weaviate / Chroma (Vector Databases for Semantic Search & RAG)

- **Role in Script Kitty:** Vector databases are indispensable for enabling semantic search, similarity retrieval, and building robust RAG (Retrieval Augmented Generation) pipelines, allowing Script Kitty's LLMs to access and synthesize information beyond their initial training data.
- **Microscopic Details:**
  - **Vector Embeddings:** Stores high-dimensional numerical representations (vectors) of text chunks, images, code snippets, or other data generated by specialized `embedding models` (e.g., `sentence-transformers`, `OpenAI embeddings`).

- **Semantic Search (k-NN):** Allows querying for items that are *semantically similar* to a given input vector, rather than just keyword matching.
- **RAG (Retrieval Augmented Generation) Pipelines:** This is critical for Script Kitty's "vast source of information" capability. When an LLM (Nexus, Core, Skills) needs factual information:
  - The query is embedded into a vector.
  - This vector is used to search the vector database for the `top-k` most semantically similar chunks of knowledge (e.g., from `Armory`'s digested web data, `Core`'s GKGS node descriptions, `Skills`' code snippets, `Guardian`'s policy documents).
  - The retrieved chunks are then inserted into the LLM's context window as supplementary information, allowing the LLM to generate more accurate, grounded, and up-to-date responses.
- **Use Cases for Script Kitty:**
  - **Core:** Stores vector embeddings of GKGS nodes/relationships for semantic search of relevant knowledge during planning and reasoning.
  - **Armory:** Stores embeddings of digested web content, research papers, API documentation, and tool manifests for semantic retrieval when `Core` needs to find relevant external information or tools.
  - **Skills:** Can store embeddings of past successful code solutions or problem-solving approaches to retrieve relevant examples for code generation.
  - **Guardian:** Stores embeddings of policy documents and safety guidelines for semantic search to ensure compliance.
- **Indexing Algorithms:** Utilizes efficient Approximate Nearest Neighbor (ANN) algorithms (e.g., HNSW, IVF) for fast, scalable similarity search over billions of vectors.
- **Data Types:** Stores `text chunks`, `vectors`, and associated `metadata` (e.g., source URL, timestamp, author).

## 5. MinIO / Ceph / AWS S3 / GCP GCS (Object Storage for Data Lake)

- **Role in Script Kitty:** Object storage forms the scalable, durable, and cost-effective foundation of Script Kitty's Data Lake. It is the central repository for all raw, processed, and derived data assets across the entire system.
- **Microscopic Details:**
  - **Blob Storage:** Stores data as "objects" (blobs) within "buckets." Objects are typically immutable and are accessed via unique keys (paths).
  - **Massive Scalability:** Designed to handle exabytes of data and trillions of objects, making it ideal for Script Kitty's ever-growing datasets (web scrapes, training data, model artifacts, logs, metrics).

- **High Durability:** Offers extremely high data durability (e.g., 99.999999999% for S3) through replication and checksumming, ensuring Script Kitty's valuable data is never lost.
- **Cost-Effectiveness:** Generally cheaper than block or file storage, especially for archival or infrequently accessed data (leveraging lifecycle policies).
- **API Compatibility (S3-compatible):** MinIO and Ceph provide S3-compatible APIs, allowing Script Kitty's components to interact with them using standard S3 client libraries, regardless of whether it's self-hosted or cloud-based.
- **Use Cases for Script Kitty:**
    - **Armory:** Stores raw web scrapes, parsed HTML, extracted text, and structured data before further processing.
    - **Skills:** Stores large model training datasets, intermediate processed data, and trained model checkpoints.
    - **Guardian:** Archives raw performance logs, audit trails, and feedback data.
    - **Foundry:** Stores Docker images (for self-hosted registries), CI/CD artifacts, DVC data versions, and `Feast` offline feature store data.
    - **Global Knowledge Graph (GKGS):** Can store large text documents referenced by GKGS nodes, with `Weaviate/Chroma` storing their embeddings.
- **Lifecycle Policies:** Configured to automatically transition data to colder storage tiers (e.g., S3 Glacier, GCS Archive) or expire it after a defined period, optimizing costs.
- **Versioning:** Object versioning ensures that every change to an object creates a new version, providing a history and enabling easy rollbacks of data (complementing `DVC`).

## 6. DVC (Data Version Control)

- **Role in Script Kitty:** DVC (Data Version Control) integrates with Git to provide version control for datasets and machine learning models, ensuring reproducibility, traceability, and collaborative development across Script Kitty's data-driven workflows.
- **Microscopic Details:**
    - **Git Integration:** DVC `links` small `.dvc` files in Git to large data/model files stored in object storage (e.g., MinIO, S3). Git tracks the lightweight `.dvc` files and their checksums, while DVC manages the actual data.
    - **Data Reproducibility:** Allows any Script Kitty developer or automated pipeline (e.g., `Guardian`'s retraining pipeline, `Foundry`'s CI/CD) to checkout a specific version of a dataset, model, or pipeline from Git, ensuring that experiments and deployments are fully reproducible.
    - **Pipeline Definition (`dvc.yaml`):** Defines reproducible data processing and ML training pipelines as a series of stages. Each stage specifies its `dependencies` (input data, code), `commands` (to run), and `outputs` (processed data, models).

This is crucial for automating `Guardian`'s continuous training loops for `Skills` models.

- ○ **Artifact Tracking:** DVC automatically tracks and versions the output artifacts of each pipeline stage, including trained models, preprocessed datasets, and evaluation metrics.
- ○ **Use Cases for Script Kitty:**
    - ■ **Skills:** Versions training datasets, validation datasets, and trained model artifacts for all specialized ML models.
    - ■ **Guardian:** Ensures that the data used for retraining specific models is versioned and reproducible, allowing `Guardian` to track the exact lineage of improved models and perform A/B testing on different data versions.
    - ■ **Armory:** Versions large datasets collected from web scraping or external APIs.
    - ■ **Core:** Can version datasets used for training its own planning heuristics or knowledge graph embeddings.

## 7. Feast (Feature Store)

- ● **Role in Script Kitty:** Feast serves as Script Kitty's centralized, standardized, and discoverable repository for machine learning features, ensuring consistency, reusability, and low-latency access for both training and inference across `Skills` and `Guardian` agents.
- ● **Microscopic Details:**
    - ○ **Feature Definitions:** Features are defined declaratively in YAML files, specifying their `name`, `data type`, `entity (e.g., user_id, tool_id)`, `source` (e.g., Kafka stream, PostgreSQL table), and any `transformation logic`.
    - ○ **Offline Store:** Uses `Foundry`'s `Object Storage` (e.g., S3) or `PostgreSQL` as the historical store for large volumes of features used during model training by `Guardian`. It ensures point-in-time correctness, meaning features used for training accurately reflect the state of the world at the time the events occurred, preventing data leakage.
    - ○ **Online Store:** Uses a low-latency key-value store like `Redis` (or Cassandra for larger scale) to serve fresh features for real-time model inference. `Skills` agents, when performing a prediction, can quickly fetch the latest feature values.
    - ○ **Consistency:** Guarantees that the features consumed by a model during training (offline) are precisely the same as those consumed during real-time inference (online), eliminating "training-serving skew," a common ML operational pitfall.
    - ○ **Feature Discovery & Reusability:** Provides a catalog of available features, preventing redundant feature engineering and promoting reuse across different Script Kitty ML models.
    - ○ **Use Cases for Script Kitty:**
        - ■ **Skills:** Serves features like "user's recent activity score," "tool's historical success rate," "code snippet complexity score" to `Skills` agents for

real-time predictions (e.g., in a recommendation model, code generation confidence).

- **Guardian:** Provides aggregated features like "model inference error rate over last 24 hours," "feedback score trends," "anomalous system call counts" to models that predict the need for retraining or detect system health degradation.
- **Core (Potential):** Could expose features related to `Core`'s internal state or past planning successes/failures to inform its own meta-learning or self-improvement.

---

**III. Compute Resources & Specialized Processors**

**1. General-Purpose CPUs (for all Microservices)**

- **Role in Script Kitty:** Standard Central Processing Units are the workhorses for all non-ML heavy computational tasks and for managing the overall distributed system. Every microservice in Script Kitty, including `Nexus`'s API gateways, `Core`'s orchestration logic, `Armory`'s web scrapers, `Skills`' sandboxed execution environments (for non-ML code), `Guardian`'s policy engine, and `Foundry`'s MLOps tools, relies on CPU power.
- **Microscopic Details:**
  - **Compute Instances:** Provided by Kubernetes worker nodes (VMs in cloud, or physical servers). These are typically multi-core CPUs with varying clock speeds and cache sizes.
  - **Workload Types:**
    - **I/O Bound:** `Nexus Channel Adapters` (handling many concurrent connections), `Armory`'s web fetching, `Foundry`'s log collectors. These require efficient context switching and network processing.
    - **CPU Bound (Light/Medium):** `Core`'s rule-based planning, `Guardian`'s policy engine (for `OPA` rules), `Nexus`'s pre-processing, `Foundry`'s CI/CD agents. These benefit from higher clock speeds and more cores.
    - **Batch Processing:** `Armory`'s data digestion, `Skills`'s data analysis (non-ML heavy), `Guardian`'s offline evaluation pipelines.
  - **Resource Requests & Limits (Kubernetes):** Each Script Kitty `Pod` is configured with `cpu.requests` and `cpu.limits`. `requests` define the minimum CPU guarantee, while `limits` cap the CPU usage to prevent any single component from starving others. This ensures fair resource allocation and stable performance across Script Kitty.
  - **Concurrency:** Many Script Kitty services are designed to be highly concurrent, utilizing asynchronous programming (e.g., Python's `asyncio`, Go's goroutines) to make efficient use of CPU cores by interleaving I/O operations with computation.

**2. GPUs (Graphics Processing Units - for LLMs & Deep Learning Models)**

- **Role in Script Kitty:** GPUs are the specialized accelerators essential for running large language models (LLMs) and other deep learning models that form the "brain" of Script Kitty's intelligent components, particularly `Nexus`, `Core`, `Skills`, and `Guardian`.
- **Microscopic Details:**
  - **Parallel Processing:** GPUs excel at massively parallel computation, executing thousands of simple arithmetic operations simultaneously. This architecture is perfectly suited for the matrix multiplications and convolutions that dominate deep learning inference and training.
  - **LLM Inference:**
    - **Nexus:** The `Nexus LLM` (Mistral-7B/Llama 3-8B) relies on GPUs for low-latency text generation and NLU (intent/entity extraction) to provide a responsive user experience.
    - **Core:** The `Core LLM` (Llama 3-70B, GPT-4o API via adapter) for high-level planning and reasoning requires significant GPU power for complex task decomposition.
    - **Skills:** Various specialized deep learning models within `Skills` (e.g., computer vision, advanced NLP, code generation) are GPU-accelerated.
    - **Guardian:** The `Ethical LLM` and models for bias detection might also leverage GPUs.
  - **Model Training:**
    - **Guardian:** Orchestrates the retraining of `Skills` models, `Nexus` LLM, and potentially `Core`'s LLM, which is a highly GPU-intensive process. `Foundry`'s `Kubeflow Pipelines` and `Ray Tune` manage these GPU-accelerated training jobs.
  - **Memory (VRAM):** GPUs come with dedicated high-bandwidth memory (VRAM). The size of the VRAM dictates the largest LLM or deep learning model that can be loaded and processed efficiently. Model quantization techniques (e.g., 4-bit, 8-bit) are used to reduce VRAM requirements, allowing larger models to fit on available GPUs.
  - **GPU Drivers & CUDA:** Proper NVIDIA GPU drivers and CUDA (Compute Unified Device Architecture) toolkit are essential for deep learning frameworks (PyTorch, TensorFlow) to communicate with and leverage the GPU hardware.
  - **Kubernetes GPU Scheduling:** Kubernetes is configured with `NVIDIA Device Plugin` (or similar for other vendors) to expose GPUs as schedulable resources. `Pods` specify `resources.limits.nvidia.com/gpu: 1` (or more) to request GPU allocation.
  - **Dedicated GPU Nodes:** `Foundry` maintains dedicated Kubernetes worker nodes equipped with multiple high-performance GPUs (e.g., NVIDIA A100s, H100s) specifically for Script Kitty's demanding ML workloads.

**3. Quantum Computing Access (Via Armory - Future/Specialized)**

- **Role in Script Kitty:** While not a core operational component today, `Armory` is designed with the foresight to access specialized quantum computing resources for problems where classical computation is intractable or highly inefficient (e.g., complex optimization problems, advanced materials science simulations, specific cryptographic challenges).
- **Microscopic Details:**
    - **Service via Armory:** `Armory` acts as the gateway. When `Core` identifies a sub-problem suitable for quantum computation (based on the `GKGS` knowledge base), it dispatches the request to `Armory`.
    - **SDKs/APIs:** `Armory` utilizes vendor-specific SDKs like `Qiskit` (IBM Quantum), `Cirq` (Google Quantum AI), or `PennyLane` (Xanadu) to programmatically interact with quantum hardware (real quantum computers or simulators).
    - **Job Submission & Management:** `Armory` manages the entire lifecycle of quantum jobs: translating classical problem definitions into quantum circuits, submitting jobs to quantum queuing systems, monitoring their execution, retrieving results, and translating quantum output back into a format usable by `Core` or `Skills`.
    - **Access Keys & Authentication:** `Armory` securely manages API keys and authentication tokens for quantum cloud services, retrieving them from `HashiCorp Vault`.
    - **Hybrid Algorithms:** Focus will be on hybrid classical-quantum algorithms (e.g., Variational Quantum Eigensolver (VQE), Quantum Approximate Optimization Algorithm (QAOA)) where the quantum computer performs the hardest part of the calculation, and a classical computer (e.g., a `Skills` agent) optimizes the overall process.
    - **Cost & Latency Considerations:** Quantum computing is inherently expensive and often has high queueing latency. `Core`'s planning engine, informed by `Armory`'s knowledge of quantum resource availability and cost, will only invoke quantum computation for problems where it offers a significant advantage.
    - **Sandboxing:** Quantum job execution, if involving custom code, might still occur within `Skills`-like sandboxed environments, orchestrated by `Armory`.

---

### IV. Open-Source Tools & Frameworks (Software Specific)

### 1. FastAPI / Flask (Python Web Frameworks)

- **Role in Script Kitty:** Lightweight, high-performance Python web frameworks used for building RESTful APIs for various microservices, particularly within `Nexus` (user-facing API), `Foundry` (internal API for MLOps orchestration), and `Guardian` (human feedback interface backend).
- **Microscopic Details:**

- **Asynchronous Support (FastAPI):** FastAPI's `async`/`await` support is crucial for building high-concurrency, non-blocking APIs, allowing `Nexus` to handle thousands of concurrent user requests (e.g., from `Channel Adapters`) efficiently without blocking the event loop.
- **Pydantic Integration (FastAPI):** FastAPI's tight integration with Pydantic for data validation and serialization ensures strict schema adherence for all incoming requests and outgoing responses across Script Kitty's internal APIs (e.g., `ScriptKittyMessageSchema`, `ToolManifestSchema`). This guarantees data integrity and reduces runtime errors.
- **Auto-generated OpenAPI/Swagger Docs (FastAPI):** Simplifies API development and consumption for other Script Kitty services and external integrations, as `FastAPI` automatically generates interactive API documentation.
- **Lightweight & Minimal (Flask):** For simpler internal services or quick prototypes where the full power of `FastAPI` isn't strictly necessary.
- **Microservice Architecture:** Both frameworks are well-suited for building small, independent microservices, aligning perfectly with Script Kitty's modular design.

## 2. Express.js (Node.js Web Framework)

- **Role in Script Kitty:** A fast, unopinionated, minimalist web framework for Node.js, primarily used in `Nexus` for building high-concurrency I/O proxies or WebSocket servers, especially for web-based chat interfaces.
- **Microscopic Details:**
  - **Non-blocking I/O:** Node.js's single-threaded, event-driven, non-blocking I/O model makes it highly efficient for handling a large number of concurrent connections (e.g., from `WebChatAdapters` using WebSockets) with low overhead.
  - **WebSocket Integration:** Express.js combined with libraries like `ws` or `socket.io` makes it straightforward to build real-time bidirectional communication channels required for interactive web chat within `Nexus`.
  - **API Gateway/Proxy:** Can be used by `Nexus` as a performant reverse proxy or API gateway to forward requests to downstream Python-based microservices, effectively handling the high volume of incoming web requests before they hit the Python stack.
  - **Middleware:** The middleware pattern allows for flexible request processing (e.g., authentication, logging, request pre-processing) before routing to the final handler.

## 3. vLLM / Ollama (LLM Inference Servers)

- **Role in Script Kitty:** Dedicated, high-performance LLM inference engines essential for serving the various Large Language Models utilized by `Nexus` (NLU & Response), `Core`

(Planning), `Skills` (Code Generation, Specialized Models), and `Guardian` (Ethical LLM). They ensure low-latency and high-throughput LLM operations.

- **Microscopic Details:**
  - **vLLM:**
    - **PagedAttention:** A groundbreaking attention algorithm that efficiently manages KV cache memory, dramatically increasing throughput for LLM serving by reducing memory fragmentation. This means Script Kitty can serve more LLM requests per GPU, reducing overall inference cost.
    - **Continuous Batching:** Dynamically batches incoming requests without waiting for a full batch, maximizing GPU utilization by keeping the GPU busy with computation.
    - **CUDA Kernels:** Optimized custom CUDA kernels for various LLM operations, leading to significantly faster inference speeds compared to generic implementations.
    - **Support for Diverse Models:** Supports a wide range of popular open-source LLMs (e.g., Llama, Mistral, CodeLlama), allowing Script Kitty to dynamically swap LLM models based on specific task requirements (e.g., a smaller, faster model for `Nexus` NLU, a larger, more powerful one for `Core`'s complex planning).
  - **Ollama:**
    - **Local/Edge Deployment:** Excellent for running smaller LLMs on local machines or edge devices, potentially for a lightweight "dev" environment of Script Kitty, or for future distributed edge components.
    - **Easy Setup:** Simplifies the process of downloading and running LLMs, abstracting away complex dependency management.
    - **API Interface:** Provides a simple API for inference, making it easy for Script Kitty components to integrate.
  - **Quantization Support:** Both support quantized models (e.g., 4-bit, 8-bit), which significantly reduce the memory footprint and computational requirements of LLMs, allowing Script Kitty to run larger models on fewer or less powerful GPUs.
  - **HTTP/gRPC API:** Provide a simple interface that Script Kitty's components (`Nexus`, `Core`, `Skills`, `Guardian`) can query to perform LLM inference, abstracting away the complexities of GPU management and model loading.

### 4. Hugging Face Transformers (ML Model Library & Tooling)

- **Role in Script Kitty:** The foundational library for accessing, fine-tuning, and deploying a vast array of pre-trained transformer models, crucial for various NLP tasks across `Nexus`, `Core`, `Armory`, `Skills`, and `Guardian`.
- **Microscopic Details:**
  - **Model Hub Integration:** Provides seamless access to the Hugging Face Model Hub, a repository of thousands of pre-trained models (e.g., `bert-base-uncased`, `roberta-base`, `bart-large-cnn`, `t5-small`,

`CodeLlama`, `Mistral`). Script Kitty leverages this for quick integration of state-of-the-art models.
- **Unified API:** Offers a consistent API for loading, fine-tuning, and running inference with diverse transformer architectures, simplifying the development of ML components.
- **Specific Use Cases in Script Kitty:**
    - **Nexus:** For the `Traditional Classifier` fallback in `NLU` (e.g., fine-tuned `bert-base-uncased` for fast intent classification).
    - **Armory:** For `Summarization` models (`bart-large-cnn`, `t5-small`) to digest web content, and for `Named Entity Recognition (NER)` or `Relation Extraction (RE)` models to structure extracted information.
    - **Skills:** For specialized `NLP models` (e.g., for sentiment analysis, text generation, question answering), and for base `Code-focused LLMs` like `CodeLlama` or `StarCoder2` for code generation.
    - **Core:** Could be used for specific `semantic embedding` models to generate vectors for `Weaviate/Chroma` from `GKGS` nodes.
    - **Guardian:** For `Bias Detection` models or `Ethical LLM` fine-tuning.
- **Tokenizers:** Provides optimized tokenizers for each model, handling text pre-processing (splitting text into tokens, converting to IDs) correctly, which is critical for accurate LLM input.
- **Pipelines API:** Simplifies common NLP tasks (e.g., `pipeline("text-generation")`, `pipeline("summarization")`), allowing for rapid prototyping and deployment of ML capabilities.
- **Fine-tuning Utilities:** Provides tools and scripts for fine-tuning pre-trained models on custom datasets. `Guardian` leverages this for continuous improvement of `Nexus` and `Skills` models using Script Kitty's experience data.

### 5. spaCy (Advanced NLP Library)

- **Role in Script Kitty:** A highly optimized NLP library for production-ready language processing, complementing `Hugging Face Transformers` by offering fast, rule-based, and custom NLP components, particularly in `Nexus` and `Armory`.
- **Microscopic Details:**
    - **Efficiency:** Designed for speed, making it suitable for high-throughput text processing in `Nexus` (pre-processing, basic entity recognition) and `Armory` (faster parsing of scraped content).
    - **Custom NER & Rule-based Matching:** Allows for the creation of custom Named Entity Recognition (NER) models or rule-based `Matcher` components. This is valuable in `Armory` for extracting specific, domain-specific entities from unstructured text that might be hard for generic LLMs to reliably capture (e.g., product codes, specific technical parameters).

- **Dependency Parsing & Part-of-Speech Tagging:** Provides structural linguistic analysis (who did what to whom), which can enrich entities extracted by `Armory` or provide deeper context for `Core`'s reasoning if specific linguistic patterns are relevant.
- **Pre-trained Models:** Offers high-quality, pre-trained statistical models for various languages.
- **Integration:** Can work in conjunction with LLMs. For instance, an LLM might identify the broad intent, and `spaCy` might then precisely extract highly specific entities using custom rules.

## 6. Playwright / Selenium (Headless Browser Automation)

- **Role in Script Kitty:** Essential tools for `Armory` to perform intelligent web scraping and data ingestion from dynamic, JavaScript-heavy websites. They simulate a real user's browser interaction.
- **Microscopic Details:**
    - **Headless Mode:** Runs web browsers (Chromium, Firefox, WebKit for Playwright; Chrome, Firefox, Edge for Selenium) without a visible GUI, making them suitable for server-side automation.
    - **JavaScript Execution:** Crucially, they execute JavaScript on web pages, allowing `Armory` to scrape data from single-page applications (SPAs) or content loaded dynamically after page load.
    - **Browser Interaction:** `Armory` can programmatically interact with web elements: clicking buttons, filling forms, scrolling, waiting for elements to load, and navigating multi-page workflows.
    - **Cross-Browser Support (Playwright):** Playwright supports all major rendering engines (Chromium, Firefox, WebKit) with a single API, ensuring `Armory`'s scraping logic works across diverse web technologies.
    - **Screenshot & PDF Generation:** Can capture screenshots or generate PDFs of web pages, useful for `Armory` to provide visual context during data digestion or for `Guardian`'s auditing purposes.
    - **Network Interception (Playwright):** Playwright can intercept network requests and responses, allowing `Armory` to filter out irrelevant resources (e.g., ads, analytics scripts), modify requests, or directly access API responses that populate the page.
    - **Session Management:** `Armory` can manage browser sessions, including cookies and local storage, to maintain login states or handle session-specific data.
    - **Resource Management:** Running headless browsers can be resource-intensive. `Armory` is designed to run these in isolated, containerized environments managed by Kubernetes, with appropriate resource limits.

## 7. BeautifulSoup4 / lxml (HTML Parsing Libraries)

- **Role in Script Kitty:** Efficient and robust Python libraries used by `Armory` to parse HTML and XML content, extract specific data elements, and navigate document structures.
- **Microscopic Details:**
  - **Parsing HTML/XML:** Take raw HTML strings (retrieved by `Playwright/Selenium` or `requests`) and convert them into a parse tree, allowing for easy navigation and searching.
  - **CSS Selectors / XPath:** Support powerful querying mechanisms (CSS selectors for BeautifulSoup4/lxml, XPath for lxml) to locate specific elements based on their tags, classes, IDs, or attributes. This is essential for `Armory` to extract structured data from web pages.
  - **Robustness:** Designed to handle malformed HTML, gracefully parsing even broken or inconsistent markup, which is common in real-world web scraping.
  - **Tree Traversal:** Allows `Armory` to navigate the document tree (e.g., `parent`, `children`, `siblings`) to find related information once an element is located.
  - **Performance (lxml):** `lxml` is written in C and is significantly faster for large HTML documents, critical for `Armory`'s high-volume data ingestion.

## 8. Trafilatura / boilerpy3 (Content Extraction Libraries)

- **Role in Script Kitty:** Specialized libraries used by `Armory` to extract the main content (e.g., articles, blog posts) from web pages, removing boilerplate (navigation, ads, footers) and formatting it cleanly.
- **Microscopic Details:**
  - **Boilerplate Removal:** Algorithms within these tools analyze the HTML structure and content density to identify and strip away non-content elements, leaving only the primary article text. This is crucial for `Armory` to focus on relevant information for `Core` and the `GKGS`.
  - **Metadata Extraction:** Can often extract valuable metadata like author, publication date, title, and images associated with the main content.
  - **Readability Enhancement:** Often re-formats the extracted text for better readability, removing excessive line breaks or unusual characters.
  - **Use Cases for Script Kitty:** Enables `Armory` to feed cleaner, higher-quality text into `Hugging Face Transformers` for summarization, NER, or into `Weaviate/Chroma` for embedding, improving the quality of Script Kitty's knowledge acquisition.

## 9. Pandas, NumPy, SciPy (Python Scientific Computing Libraries)

- **Role in Script Kitty:** The bedrock of numerical computing and data manipulation in Python, extensively used by `Skills` agents for data analysis, transformation, and statistical operations, and by `Guardian` for evaluation metrics.
- **Microscopic Details:**

- **Pandas:**
  - **DataFrame:** The primary data structure, a tabular, column-oriented data container. `Skills` agents use DataFrames for cleaning, transforming, aggregating, and analyzing structured and semi-structured data (e.g., from `Armory`'s extractions, or `Feast` features).
  - **Data Cleaning & Preprocessing:** Provides powerful tools for handling missing values, filtering rows/columns, merging datasets, and reshaping data – all critical steps before feeding data into ML models or performing analysis.
  - **Time Series Analysis:** Strong support for time-series data, useful for `Guardian` to analyze historical performance metrics or `Skills` for time-dependent forecasting.
- **NumPy:**
  - **ndarray (N-dimensional array):** The fundamental data structure for numerical computation. All deep learning frameworks (PyTorch, TensorFlow) are built on top of NumPy arrays. `Skills` and `Guardian` use NumPy for efficient vectorized operations, mathematical functions, and linear algebra.
  - **Performance:** Highly optimized C/Fortran backend ensures very fast numerical computations, crucial for processing large datasets.
- **SciPy:**
  - **Scientific Computing Modules:** Provides modules for optimization, linear algebra, integration, interpolation, special functions, FFT, signal and image processing, and other scientific and engineering tasks.
  - **Statistical Analysis:** `Guardian` leverages SciPy for advanced statistical tests, hypothesis testing, and probability distributions when evaluating model performance or detecting subtle anomalies. `Skills` might use it for specialized statistical modeling.

## 10. Dask / Apache Spark (PySpark) (Distributed Data Processing)

- **Role in Script Kitty:** Frameworks for distributed data processing, enabling `Skills` and `Guardian` to handle datasets that exceed the memory capacity of a single machine, or require accelerated processing across a cluster of `Foundry`'s `Kubernetes` worker nodes.
- **Microscopic Details:**
  - **Dask:**
    - **Python-Native:** Integrates seamlessly with `NumPy`, `Pandas`, and `scikit-learn`, making it easier for Python-centric `Skills` agents to scale their workflows.
    - **Lazy Evaluation:** Operations are built into a task graph and executed only when results are needed, allowing for optimization.

- **Flexible Deployment:** Can run on a single machine, a cluster of VMs/containers (managed by Kubernetes), or cloud providers. `Foundry` would orchestrate Dask clusters on `Kubernetes` for specific `Skills` tasks.
- **Use Cases:** Distributed `Pandas DataFrames`, `NumPy Arrays`, and custom task graphs for `Skills` agents performing large-scale data cleaning, feature engineering, or statistical analysis.
  - **Apache Spark (PySpark):**
    - **Unified Analytics Engine:** Provides APIs for SQL, streaming, MLlib (machine learning), and GraphX (graph processing).
    - **RDDs (Resilient Distributed Datasets) / DataFrames:** Core data structures for distributed computation.
    - **In-Memory Processing:** Optimizes performance by keeping data in memory across iterations.
    - **Fault Tolerance:** Recovers from node failures gracefully.
    - **Use Cases:** `Guardian` for large-scale data ingestion and batch processing of logs and metrics from the `Data Lake` for anomaly detection or training `Foundry`'s internal models. `Skills` could use it for complex ETL (Extract, Transform, Load) operations on very large datasets before feeding them to ML models.
  - **Kubernetes Integration:** `Foundry` would provision and manage Dask or Spark clusters as Kubernetes custom resources or through operators (e.g., `Spark Operator`), dynamically spinning up resources as needed by `Skills` or `Guardian`.

## 11. Matplotlib, Seaborn, Plotly (Data Visualization Libraries)

- **Role in Script Kitty:** Essential tools for `Skills` agents to generate insightful visualizations of data analysis results, and for `Guardian` to present performance metrics, evaluation results, and detected biases in human-readable formats (e.g., within `Guardian`'s `Human Oversight & Feedback Interface`).
- **Microscopic Details:**
  - **Matplotlib:** The foundational plotting library in Python, offering granular control over every aspect of a plot. `Skills` uses it for static, high-quality plots for reports.
  - **Seaborn:** Built on Matplotlib, providing a higher-level interface for creating aesthetically pleasing statistical graphics (e.g., heatmaps, violin plots, regression plots). Useful for `Skills` to quickly explore relationships in data or for `Guardian` to visualize model confidence distributions.
  - **Plotly:**
    - **Interactive Visualizations:** Crucial for the `Guardian Human Oversight & Feedback Interface`. Plotly generates interactive

plots (zoom, pan, hover tooltips) that can be embedded directly into web applications (e.g., using `Dash` or `Streamlit`). This allows human operators to drill down into Script Kitty's performance or review detected anomalies.

- **Dashboarding:** `Plotly Dash` enables building custom interactive dashboards for `Guardian`'s monitoring, providing a more dynamic and actionable view than static images.
○ **Generative Task Support:** `Skills` agents, when tasked with generating reports or analytical summaries, can programmatically generate these visualizations as part of their output.

## 12. Pydantic (Data Validation & Settings Management)

- **Role in Script Kitty:** A Python library for data validation and settings management using Python type hints, ensuring data integrity and consistency across Script Kitty's numerous internal APIs and data flows.
- **Microscopic Details:**
  ○ **Type Hinting for Data Structures:** Defines data models (e.g., `ScriptKittyMessageSchema`, `TaskRequestSchema`, `ToolManifestSchema`) using standard Python type hints.
  ○ **Automatic Validation:** Automatically validates input data against the defined schemas, raising informative errors if data does not conform. This reduces boilerplate validation code and makes APIs more robust.
  ○ **Data Parsing & Serialization:** Automatically parses input data (e.g., from JSON strings) into Python objects and can serialize Python objects back into JSON or other formats.
  ○ **FastAPI Integration:** Tightly integrated with `FastAPI`, making it the default for request/response validation.
  ○ **Use Cases Across Script Kitty:**
    - **Nexus:** Validates incoming user messages and outgoing responses.
    - **Core:** Validates `TaskRequests` from `Nexus` and `TaskResults/Updates` from `Armory`/`Skills`/`Guardian`.
    - **Armory:** Validates `ToolManifests` and extracted structured data from web scraping.
    - **Skills:** Validates input parameters for specialized models and the structure of generated code.
    - **Guardian:** Validates `PolicyRules` and `Feedback` data.
    - **Foundry:** Validates configuration files and internal API requests.
  ○ **Configuration Management:** Can also be used to define application settings, reading values from environment variables or `.env` files, simplifying deployment configurations across Kubernetes environments.

## 13. LangChain / AutoGen (Agent Orchestration Frameworks)

- **Role in Script Kitty:** Frameworks that provide abstractions and tools for building and orchestrating complex multi-agent workflows, particularly within `Script Kitty.Core` for its goal decomposition and agent orchestration engine.
- **Microscopic Details:**
  - **Agent Abstraction:** Provides a way to define and interact with "agents" (which map to Script Kitty's specialized components like `Armory`, `Skills`, `Guardian`) with defined capabilities and tool usage.
  - **Chains/Graphs (LangChain):** Allows `Core` to define sequences or graphs of operations where the output of one agent or tool becomes the input for another, enabling multi-step problem-solving.
  - **Conversational Memory Integration:** Simplifies the integration of conversational memory (from `Nexus`'s `SessionManager`) into LLM-driven agents.
  - **Tool Calling:** Provides mechanisms for LLMs to "call" external tools (e.g., `Armory`'s web scraping, `Skills`' code execution) based on their understanding of the user's intent. `Core` leverages this to dynamically invoke the correct specialized agents.
  - **AutoGen (Multi-Agent Conversation):** Focuses on enabling complex, multi-agent conversations and collaborations. `Core` could use AutoGen to:
    - Orchestrate a debate between a `Skills` agent (proposing a solution) and a `Guardian` agent (evaluating ethical implications).
    - Have `Armory` (researcher agent) consult `Skills` (coder agent) to refine a data extraction script.
    - This directly facilitates the "group chat" illusion and the internal collaborative problem-solving nature of Script Kitty.
  - **Prompt Engineering Utilities:** Offers tools for managing and generating complex prompts for LLMs, including few-shot examples and structured output instructions.
  - **Observability Integration:** Often includes callbacks or hooks for integrating with observability tools (like `LangSmith` for LangChain, or custom integrations with `OpenTelemetry` for `Foundry`) to trace agent interactions.
  - **Dynamic Workflow Generation:** While `Core` has its own HTN planner, these frameworks provide a lower-level API to programmatically construct and execute agent workflows, allowing `Core`'s LLM to dynamically generate execution plans.

### 14. ply (Python Lex-Yacc - for Symbolic Planning Fallback)

- **Role in Script Kitty:** A Python implementation of `lex` and `yacc`, used by `Script Kitty.Core` to parse simple domain-specific languages, specifically for its symbolic HTN (Hierarchical Task Network) planning fallback mechanism or for custom rule syntax.
- **Microscopic Details:**

- ○ **Lexical Analysis (Lexer):** Breaks down input text into a sequence of tokens (e.g., keywords, identifiers, operators). `Core` would use this to tokenize a simplified PDDL-like (Planning Domain Definition Language) plan or a custom rule syntax.
- ○ **Syntactic Analysis (Parser):** Builds a parse tree (Abstract Syntax Tree - AST) from the token stream, ensuring the input conforms to a predefined grammar. This allows `Core` to formally validate and interpret symbolic plans generated by its `HTN Planner` or to parse domain-specific rules.
- ○ **Use Cases in Script Kitty:**
    - ■ **Core's HTN Planner:** If the `Core LLM` generates a plan in a structured, symbolic format, `ply` can be used to parse and validate that plan against known planning domains defined within the `GKGS`.
    - ■ **Custom Rule Engines:** For small, performance-critical rule sets that `Core` might use as guardrails or specific heuristics, `ply` can parse custom rule syntax.
- ○ **Formality & Verifiability:** Provides a way to introduce formal, verifiable components into `Core`'s planning logic, offering a robust fallback or cross-check for LLM-generated plans, especially for safety-critical operations.

### 15. Kubernetes Python Client (Python SDK for Kubernetes API)

- ● **Role in Script Kitty:** The official Python client library for interacting with the Kubernetes API, used by `Foundry` to programmatically manage and orchestrate all Script Kitty components and their underlying resources.
- ● **Microscopic Details:**
    - ○ **API Interaction:** Allows `Foundry` services (e.g., `Agent Orchestration & Lifecycle Management` in `Core`, `Resource Provisioning` in `Armory`, `CI/CD` components) to directly call Kubernetes API endpoints.
    - ○ **Resource Management:**
        - ■ `Foundry` can dynamically `create`, `delete`, `update`, `get`, and `watch` Kubernetes resources:
            - ■ `Deployments` for stateless services (Nexus, Core, Guardian components).
            - ■ `Jobs` for one-off tasks (e.g., `Skills` sandboxed code execution, `Guardian`'s evaluation pipelines).
            - ■ `DaemonSets` for cluster-wide agents (e.g., `Fluent Bit` for logging).
            - ■ `CustomResourceDefinitions` (CRDs) for managing Script Kitty's custom resources.
        - ■ `Armory` uses it for `Resource Provisioning` (e.g., requesting an ephemeral GPU-enabled pod).

- **Core**'s `Agent Orchestration` uses it to launch and manage `Skills` and `Armory` agents as Kubernetes `Jobs` or `Deployments`.
  - **In-Cluster Configuration:** Automatically detects Kubernetes client configuration when running inside a cluster, simplifying deployment for `Foundry`'s components.
  - **Watches & Events:** Can `watch` for changes to Kubernetes resources, allowing `Foundry` to react in real-time to changes in the cluster state (e.g., a `Pod` failing, a `Job` completing).
  - **Pythonic Interface:** Provides a Pythonic interface to the Kubernetes API, simplifying automation scripts and controllers.

## 16. Argo Workflows / Kubeflow Pipelines (Workflow Orchestration)

- **Role in Script Kitty:** Powerful workflow engines running natively on Kubernetes, crucial for orchestrating complex, multi-step, and often long-running pipelines within Script Kitty, especially for `Foundry` (MLOps), `Guardian` (training/evaluation), and `Core` (complex task execution).
- **Microscopic Details:**
  - **DAG (Directed Acyclic Graph) Definition:** Workflows are defined as DAGs, specifying a sequence of `steps` or `tasks` and their `dependencies`. Each `step` typically runs as a separate Kubernetes `Pod`.
  - **Kubeflow Pipelines:** A platform for building and deploying portable, scalable ML workflows.
    - **Component Reusability:** Defines reusable components for common ML tasks (e.g., data preprocessing, training, evaluation, model serving).
    - **Tracking & Visualization:** Provides a UI for visualizing workflow execution, parameters, and outputs, which is vital for `Guardian` to monitor training runs and `Foundry` to debug MLOps pipelines.
    - **Metadata Tracking:** Automatically tracks metadata about pipeline runs (e.g., input data, model versions, metrics), integrating with `MLflow`.
  - **Argo Workflows:** A more general-purpose workflow engine for Kubernetes.
    - **Flexible Task Types:** Supports various task types including shell commands, Docker images, and Kubernetes manifests.
    - **Dynamic Workflows:** Can generate dynamic workflows based on logic within a step.
    - **Use Cases for Script Kitty:**
      - **Guardian:** Orchestrates complex `ML training pipelines` (data ingestion -> preprocessing -> model training -> evaluation -> model registration) for `Skills` models, `Nexus` LLM, and `Core`'s internal models. Also orchestrates `evaluation pipelines` for `Bias Detection`.

- **Foundry:** Manages `CI/CD pipelines` (if not using native GitLab/GitHub actions), `data processing pipelines` (e.g., large-scale data cleansing for `Armory`), and general automation.
- **Core (Advanced):** For very complex, multi-agent tasks that span multiple `Skills` and `Armory` invocations, `Core` could dynamically generate and submit an `Argo Workflow` to `Foundry` for execution.
  - **Error Handling & Retries:** Built-in mechanisms for retrying failed steps, conditional execution, and error notifications.
  - **Scalability:** Each step runs as a separate Kubernetes Pod, allowing for parallel execution and efficient resource utilization across the cluster.

### 17. Apache Kafka / Apache Pulsar (Distributed Message Brokers)

- **Role in Script Kitty:** High-throughput, fault-tolerant, and scalable distributed streaming platforms that serve as the central asynchronous communication bus for Script Kitty, decoupling services and enabling real-time data flows.
- **Microscopic Details:**
  - **Publish-Subscribe Model:** Producers (e.g., `Nexus` for `user_activity` events, `Skills` for `task_completion` events, `Guardian` for `policy_violation` events) publish messages to topics, and consumers (e.g., `Core` for `task_results`, `Foundry` for `metrics/logs`, `Guardian` for `feedback_data`) subscribe to relevant topics.
  - **Durability:** Messages are persisted to disk, ensuring data is not lost even if consumers are offline or services crash. This is vital for audit trails and recovery.
  - **Scalability:** Horizontally scalable by adding more brokers to a cluster, handling massive volumes of events (e.g., all `Script Kitty`'s internal logs, metrics, trace spans, and agent updates).
  - **Decoupling:** Services communicate asynchronously without direct dependencies, improving resilience and allowing independent development and deployment of Script Kitty's components.
  - **Event Sourcing:** Can serve as an event store, capturing the entire history of state changes within Script Kitty (e.g., every `TaskRequest`, `TaskUpdate`, `TaskResult`).
  - **Use Cases for Script Kitty:**
    - **Internal Communication Hub (Core & Foundry):** The primary bus for asynchronous messages between `Core` and its specialized agents (`Armory`, `Skills`, `Guardian`). While `gRPC` handles synchronous request/response, `Kafka/Pulsar` handles ongoing status updates, large data streams, and event notifications.

- **Observability (Foundry):** `Fluent Bit/Fluentd` can send aggregated logs and metrics to Kafka topics before they are consumed by `Loki` or `Prometheus` for more robust ingestion.
- **Feedback & Learning Integration (Guardian):** `Guardian` ingests success/failure signals, human corrections, and new observations via Kafka topics for real-time processing.
- **Data Ingestion (Armory):** Large streams of web scraped data or external API data can be pushed to Kafka topics for consumption by `Skills` for processing.
  - **Kafka Streams / Flink / Spark Streaming:** Used for real-time processing of data from Kafka topics (e.g., by `Guardian` for real-time anomaly detection on metrics streams).

## 18. gRPC (High-Performance RPC Framework)

- **Role in Script Kitty:** The primary framework for synchronous, high-performance, and efficient inter-service communication within Script Kitty, particularly between `Nexus` and `Core`, and between `Core` and its immediate child agents (`Armory`, `Skills`, `Guardian`).
- **Microscopic Details:**
  - **Protocol Buffers (Protobuf) Integration:** `gRPC` uses `Protobuf` as its Interface Definition Language (IDL) and underlying message interchange format. This ensures strict schema definition, type safety, and very efficient binary serialization/deserialization of messages.
  - **Performance:** Built on HTTP/2, enabling features like multiplexing (multiple concurrent requests over a single TCP connection), header compression, and server push, leading to significant performance gains over traditional REST+JSON.
  - **Bidirectional Streaming:** Supports four types of service methods:
    - **Unary RPC:** (client sends one request, gets one response) - e.g., `Nexus` sends `TaskRequest` to `Core`.
    - **Server-Side Streaming:** (client sends one request, gets a stream of responses) - e.g., `Core` streams `TaskUpdate` messages to `Nexus` for a long-running task.
    - **Client-Side Streaming:** (client sends a stream of requests, gets one response) - e.g., `Armory` streams large data chunks to `Skills` for processing.
    - **Bidirectional Streaming:** (both client and server send streams) - e.g., a real-time negotiation between `Core` and `Guardian` on policy enforcement.
  - **Code Generation:** `Protobuf` compiles `.proto` definitions into code for various languages (Python, Go, Node.js), providing strongly typed client and server

stubs, which reduces development effort and eliminates common API integration errors.
- **Reliability:** The underlying HTTP/2 ensures reliable connections and efficient handling of network issues.
- **Use Cases:**
  - `Nexus` to `Core`: For high-level `TaskRequest` submission and receiving final `TaskResult`.
  - `Core` to `Armory`/`Skills`/`Guardian`: For dispatching specific sub-tasks (e.g., `ExecuteCode`, `FetchTool`, `EvaluatePolicy`) and receiving their immediate results.
  - `Foundry` internal services: For high-performance communication between MLOps components.

**19. Protocol Buffers (Protobuf - Data Serialization Format)**

- **Role in Script Kitty:** The definitive language for defining the precise schema of all structured data exchanged between Script Kitty's microservices, ensuring strict type-checking, backward/forward compatibility, and efficient binary serialization/deserialization. It is the language of communication for gRPC.
- **Microscopic Details:**
  - **Language-Neutral, Platform-Neutral:** Once a `.proto` message is defined, `Protobuf` compilers generate code in multiple programming languages (Python, Go, Java, C++, etc.), allowing different Script Kitty services to communicate seamlessly regardless of their implementation language.
  - **Schema Definition:** `.proto` files define the structure of messages (e.g., `message ScriptKittyMessage { string user_id = 1; string raw_text = 2; ... }`). This enforces a contract between services.
  - **Efficient Binary Format:** Serializes data into a compact binary format, significantly smaller and faster to parse than JSON or XML, crucial for high-throughput communication within Script Kitty.
  - **Backward/Forward Compatibility:** Supports adding new fields to messages without breaking existing services, making it easier to evolve Script Kitty's internal APIs without requiring a "big bang" upgrade.
  - **Use Cases:**
    - `ScriptKittyMessageSchema`: Standardized format for all user inputs and outputs.
    - `TaskRequestSchema`, `TaskUpdateSchema`, `TaskResultSchema`: Standardized schema for all task orchestration.
    - `ToolManifestSchema`: Precise definition of tool capabilities in `Armory`.
    - `PolicyRuleSchema`: Structure for `Guardian`'s policy definitions.
    - All internal APIs: Defines request and response payloads.

- ○ **Validation:** While `Protobuf` itself doesn't do runtime validation beyond type, combined with `Pydantic` (for Python), it ensures robust data integrity.

### 20. Prometheus (Monitoring & Alerting System)

- **Role in Script Kitty:** Prometheus is the backbone of Script Kitty's real-time metrics collection and alerting system, providing granular insights into the operational health and performance of every component.
- **Microscopic Details:**
  - ○ **Time-Series Database:** Stores metrics as time-series data (values over time).
  - ○ **Pull Model:** Prometheus typically "pulls" metrics by scraping HTTP endpoints (often `/metrics`) exposed by Script Kitty's microservices at regular intervals.
  - ○ **Exporters:** Each Script Kitty microservice incorporates a Prometheus client library to expose its internal metrics (e.g., request counts, latency histograms, error rates, queue sizes, CPU/memory usage, LLM inference times, task success rates).
  - ○ **PromQL (Prometheus Query Language):** A powerful, flexible query language for querying and aggregating time-series data. `Foundry`'s operators and `Guardian`'s monitoring components use PromQL to analyze trends, identify deviations, and define alert conditions.
  - ○ **Service Discovery (Kubernetes):** The `Prometheus Operator` (deployed by `Foundry`) automatically discovers new Script Kitty services in Kubernetes based on labels or annotations and configures Prometheus to scrape them.
  - ○ **Alertmanager:** Integrates with Prometheus to manage and route alerts. `Foundry` configures Alertmanager to send critical alerts (e.g., `Core_Planner_Errors_High`, `Skills_Sandbox_Breach`, `Nexus_Latency_Spike`) to `PagerDuty`, `Slack`, or other notification channels, enabling immediate human intervention.

### 21. Grafana (Visualization & Dashboarding Tool)

- **Role in Script Kitty:** A powerful and flexible open-source platform for data visualization and dashboarding, used by `Foundry` to provide real-time operational visibility into Script Kitty's health and by `Guardian` for presenting system performance and evaluation results.
- **Microscopic Details:**
  - ○ **Data Source Integration:** Connects to various data sources, primarily `Prometheus` (for metrics) and `Loki` (for logs), but also `PostgreSQL` or `Elasticsearch` for other analytical data.
  - ○ **Interactive Dashboards:** Provides rich, interactive dashboards with panels (graphs, charts, heatmaps, tables, single stats) that visualize Script Kitty's key performance indicators (KPIs) in real-time. Examples:

- **Overall System Health Dashboard:** Aggregated latency, error rates, resource utilization across all services.
- **Nexus Performance:** Active user sessions, NLU success rates, LLM inference times.
- **Core Planning Metrics:** Plan generation time, agent invocation frequency, planning success rate.
- **Skills Execution Dashboard:** Task completion rates, sandbox resource utilization, model-specific metrics.
- **Guardian Oversight:** Policy violation counts, model bias trends, training pipeline status.
  - ○ **Templating:** Allows creation of dynamic dashboards where users can select specific services, agents, or time ranges, providing flexible exploration of Script Kitty's operational data.
  - ○ **Alerting Integration:** Integrates with `Prometheus Alertmanager` to display triggered alerts directly on dashboards and provide visual context for incidents.
  - ○ **Query Editor:** Allows users to write `PromQL` or `LogQL` queries directly in the UI to explore data.

## 22. Fluent Bit / Fluentd (Log Collection Agents)

- **Role in Script Kitty:** Lightweight and highly efficient log processors and forwarders, deployed as DaemonSets on every Kubernetes node by `Foundry` to collect all container and node-level logs from Script Kitty's components.
- **Microscopic Details:**
  - ○ **DaemonSet Deployment:** Runs as a single pod on each Kubernetes worker node, ensuring comprehensive log collection from all running Script Kitty containers.
  - ○ **Container Log Collection:** Automatically collects logs from `/var/log/containers/*`, parsing them into structured JSON or semi-structured formats.
  - ○ **Metadata Enrichment:** Automatically enriches logs with Kubernetes metadata (pod name, namespace, container ID, image name, labels), and Script Kitty-specific tags (e.g., `script_kitty_session_id`, `script_kitty_task_id`), crucial for filtering and analysis.
  - ○ **Buffering & Retries:** Buffers logs in memory or on disk and retries sending them in case of network issues or downstream collector unavailability, preventing log loss.
  - ○ **Output Plugins:** Supports various output destinations (e.g., `Loki`, `Elasticsearch`, `Kafka`, `S3`). `Foundry` configures them to send logs to `Loki` for immediate analysis and potentially `Kafka` for long-term archival or stream processing by `Guardian`.
  - ○ **Filtering & Tagging:** Can apply filters to discard irrelevant logs or tag logs with specific labels for routing to different destinations.

### 23. Loki (Log Aggregation System for Semi-Structured Logs)

- **Role in Script Kitty:** A horizontally scalable, highly available log aggregation system optimized for storing and querying semi-structured logs, preferred by `Foundry` for Script Kitty's primary log analytics due to its Prometheus-inspired `LogQL` and cost-efficiency.
- **Microscopic Details:**
  - **Index-less Design:** Unlike traditional log systems that index full log content, Loki indexes only metadata (labels) for logs. This makes it very cost-effective for storing large volumes of logs.
  - **LogQL:** A powerful query language similar to `PromQL`. `Foundry`'s SREs and `Guardian`'s analysts use `LogQL` to filter logs by labels (e.g., `job="nexus-nlu-service"`, `level="error"`, `session_id="abc-123"`) and then perform text searches on the filtered results.
  - **Integration with Grafana:** Seamlessly integrates with `Grafana`, allowing for drilling down from `Prometheus` metrics to relevant `Loki` logs in the same dashboard, crucial for rapid root cause analysis of Script Kitty's operational issues.
  - **Storage:** Stores logs in an object storage backend (e.g., `S3`, `MinIO`), ensuring scalability and durability.
  - **Multi-tenancy:** Can support multiple "tenants" (e.g., different Script Kitty environments) with logical separation of logs.

### 24. OpenTelemetry (Observability Instrumentation & Standard)

- **Role in Script Kitty:** A vendor-neutral set of APIs, SDKs, and tools used to instrument Script Kitty's services for generating and exporting telemetry data (metrics, logs, traces), providing standardized observability across the entire complex system.
- **Microscopic Details:**
  - **Unified Telemetry:** Provides a single standard for collecting all three pillars of observability:
    - **Traces:** End-to-end request flow.
    - **Metrics:** Aggregated numerical data.
    - **Logs:** Event-based text messages.
  - **SDKs & Instrumentation Libraries:** `OpenTelemetry` SDKs are embedded within the code of every Script Kitty microservice (Python, Go, Node.js). These SDKs automatically or manually instrument:
    - **Incoming/Outgoing Requests:** Automatically generates spans for HTTP/gRPC calls, database queries, and inter-service communication.
    - **Custom Spans:** Developers can manually create custom spans around critical code blocks (e.g., `Core`'s `plan_generation_logic`, `Skills`'s `model_inference_step`) to capture granular performance data.
  - **Context Propagation:** Automatically propagates `trace_id` and `span_id` across service boundaries (e.g., via HTTP headers, gRPC metadata), linking all

related operations into a single distributed trace. This is crucial for understanding how a user query in `Nexus` translates into complex interactions across `Core`, `Armory`, `Skills`, and `Guardian`.
- ○ **Exporters:** Configured to export telemetry data to `OpenTelemetry Collectors`, which then forward them to `Prometheus`, `Loki`, and `Jaeger/Zipkin`.
- ○ **Vendor Neutrality:** Ensures that Script Kitty's observability data is not locked into any single vendor's ecosystem, providing flexibility to switch backends if needed.

### 25. Jaeger / Zipkin (Distributed Tracing Backends)

- **Role in Script Kitty:** Backend systems for storing, indexing, and visualizing distributed traces, used by `Foundry` to enable in-depth debugging and performance analysis of complex, multi-service interactions within Script Kitty.
- **Microscopic Details:**
  - ○ **Trace Storage & Querying:** Receive trace data from `OpenTelemetry Collectors` and store it in a backend (e.g., Cassandra, Elasticsearch). They provide APIs for querying traces by `trace_id`, service name, operation name, or tags.
  - ○ **UI Visualization:** Offers a web-based user interface that visualizes distributed traces as flame graphs or Gantt charts. This allows `Foundry`'s SREs to:
    - ■ Quickly identify latency bottlenecks: See which service or operation took the longest within a complex Script Kitty interaction.
    - ■ Debug distributed errors: Pinpoint exactly where an error occurred across multiple services.
    - ■ Understand complex flows: Visualize the execution path of a `Core`'s `TaskRequest` as it orchestrates calls to `Armory` and `Skills`.
  - ○ **Span Details:** Each span in the visualization provides detailed information (duration, logs, tags, process info), aiding in debugging.
  - ○ **Service Graph:** Can generate a service dependency graph based on observed traces, showing how Script Kitty's services interact.

### 26. MLflow (ML Experiment Tracking & Model Management)

- **Role in Script Kitty:** An open-source platform for managing the entire machine learning lifecycle, primarily used by `Guardian` for tracking model training experiments and by `Foundry` for managing trained model versions and registering them for deployment.
- **Microscopic Details:**
  - ○ **MLflow Tracking:**
    - ■ **Experiment Logging:** `Guardian`'s training pipelines (orchestrated by `Kubeflow Pipelines`/`Argo Workflows`) automatically log parameters (e.g., learning rate, number of epochs), metrics (e.g.,

accuracy, loss, F1-score, bias metrics), and artifacts (e.g., trained model files, evaluation plots) for each training run.

- **Reproducibility:** Records the code version, environment, and dependencies used for each run, ensuring that `Guardian` can reproduce past model training results.
- **Comparison UI:** Provides a web UI to compare different training runs side-by-side, allowing `Guardian` to assess the effectiveness of various hyperparameter choices or training data versions.

- ○ **MLflow Projects:** Defines machine learning code as reusable projects, specifying dependencies and entry points. This standardizes how `Guardian`'s training code is packaged and run.
- ○ **MLflow Models:**
  - **Model Packaging:** Provides a standard format for packaging trained models from various ML libraries (PyTorch, TensorFlow, scikit-learn).
  - **Model Registry:** `Foundry` uses the MLflow Model Registry as a central hub to manage the lifecycle of trained models:
    - **Version Control:** Registers different versions of a model (e.g., `Nexus_NLU_LLM_v1`, `Nexus_NLU_LLM_v2`).
    - **Stage Management:** Transitions models through stages (e.g., `Staging`, `Production`, `Archived`), which `Foundry`'s deployment pipelines (`Argo CD`) can use to automate deployments.
    - **Annotation:** Stores metadata about each model version (e.g., responsible team, training data source, evaluation results, approval status).
  - **Model Serving:** Supports various serving options, although `vLLM` is used for LLMs, `MLflow` helps manage the model artifacts that `vLLM` loads.
- ○ **Use Cases:** `Guardian` uses `MLflow` for the continuous evaluation and retraining loop, ensuring `Script Kitty`'s ML models are always optimized and biases are monitored.

## 27. Ray Tune / Optuna (Hyperparameter Optimization Libraries)

- **Role in Script Kitty:** Frameworks for automating the process of finding the optimal hyperparameters for Script Kitty's machine learning models, used by `Guardian` during its continuous retraining cycles.
- **Microscopic Details:**
  - ○ **Automated Search:** Automates the search for optimal hyperparameters (e.g., learning rate, batch size, number of layers, regularization strength) for `Skills` models, `Nexus` LLM fine-tuning, or `Core`'s internal models.
  - ○ **Search Algorithms:** Implement various search algorithms:
    - **Grid Search:** Exhaustively tries all combinations (often too slow).

- **Random Search:** Randomly samples combinations (more efficient than grid for high-dimensional spaces).
- **Bayesian Optimization (Optuna, Ray Tune):** Builds a probabilistic model of the objective function to intelligently select promising hyperparameters to evaluate next, significantly speeding up optimization.
- **Evolutionary Algorithms (Ray Tune):** Inspired by biological evolution.
- **Early Stopping:** Can intelligently stop unpromising training runs early based on validation metrics, saving computational resources (especially GPUs). This is critical for `Guardian` managing large-scale retraining.
- **Distributed Execution:** `Ray Tune` leverages the `Ray` distributed computing framework, allowing `Guardian` to run hyperparameter tuning jobs across `Foundry`'s `Kubernetes` cluster, parallelizing evaluations of different hyperparameter combinations.
- **Integration:** Integrates with `MLflow` for logging and tracking hyperparameter tuning experiments, and with `Kubeflow Pipelines`/`Argo Workflows` for orchestrating the tuning process.

## 28. Open Policy Agent (OPA - Policy Enforcement Engine)

- **Role in Script Kitty:** `OPA` is a lightweight, general-purpose policy engine used by `Guardian` to enforce ethical, legal, and safety controls across Script Kitty's proposed actions. It provides a declarative, auditable, and centralized way to manage system-wide policies.
- **Microscopic Details:**
  - **Rego Language:** Policies are written in `Rego`, a high-level, declarative policy language. `Rego` is optimized for evaluating complex structured data.
  - **Decision as a Service:** `Guardian`'s `Policy Enforcement Engine` exposes `OPA` as a microservice. Before `Core` executes a potentially sensitive action (e.g., accessing PII, deploying code to production, making external API calls), it sends a JSON `input` (containing the proposed action, context, user ID, etc.) to `OPA`.
  - **Policy Evaluation:** `OPA` evaluates the `input` against the loaded `Rego` policies and returns a `decision` (e.g., `allow` or `deny`, along with explanations or additional metadata).
  - **Decoupled Policy:** Policies are decoupled from application code, allowing `Guardian` to update, audit, and manage policies independently without redeploying other Script Kitty components.
  - **Use Cases in Script Kitty (Guardian):**
    - **Safety Policies:** "Deny `Skills` agent execution if code attempts to delete system files."
    - **Ethical Guidelines:** "Flag any response generated by `Nexus` that contains hate speech."

- **Access Control:** "Only `Core`'s planning engine can invoke `Armory`'s resource provisioning for GPUs."
- **Resource Usage Policies:** "Prevent `Skills` agents from consuming more than X amount of GPU time without explicit `Guardian` approval."
- **Data Governance:** "Deny access to sensitive data (e.g., PII) if the requesting agent does not have explicit `user_consent` or `security_clearance`."
  - **Bundle Service:** `OPA` can fetch policy bundles (Rego code and data) from a remote server (`Foundry`'s `Object Storage` or an internal GitOps repository), allowing for dynamic policy updates without restarting `OPA`.
  - **Auditing:** `OPA` can log all policy decisions and their inputs, providing a comprehensive audit trail for `Guardian` and compliance purposes.

## 29. Aequitas, Fairlearn, Google's What-If Tool, IBM's AI Fairness 360 (AIF360) (Bias Detection Libraries)

- **Role in Script Kitty:** These open-source toolkits are integrated into `Guardian`'s `Automated Evaluation & Training Trigger` to systematically detect, measure, and analyze biases in Script Kitty's training data and the outputs of its machine learning models (especially `Nexus`'s responses, `Skills`' code generation, and `Core`'s planning decisions).
- **Microscopic Details:**
  - **Fairness Metrics:** Provide a suite of fairness metrics (e.g., demographic parity, equalized odds, statistical parity difference, disparate impact) to quantify bias across different sensitive attributes (gender, race, age, etc.). `Guardian` uses these to identify potential disparities.
  - **Bias Identification:** Help identify groups that are unfairly disadvantaged or privileged by a model's predictions or decisions.
  - **Visualization:** Offer interactive visualizations (e.g., in `What-If Tool`, `Fairlearn`) to explore bias, allowing human operators within `Guardian` to understand the nature and extent of the problem.
  - **Mitigation Strategies:** While these tools primarily focus on detection, they often suggest or integrate with methods for bias mitigation (e.g., re-weighing training data, adversarial debiasing, post-processing predictions). `Guardian` would leverage these insights to adjust training data or fine-tuning approaches.
  - **Use Cases in Script Kitty (Guardian):**
    - **Data Analysis:** Scan `Armory`'s collected datasets or `Skills`' training data for inherent biases in representation.
    - **Model Output Evaluation:** Analyze the responses generated by `Nexus`'s LLM or the code generated by `Skills` for discriminatory language, stereotypes, or unfair outcomes across different demographics.

- **Decision Bias:** Assess `Core`'s planning decisions or `Guardian`'s own policy enforcements for unintended biases.
  - **Integration into Pipelines:** These tools are integrated into `Guardian`'s `Kubeflow Pipelines`/`Argo Workflows` as dedicated evaluation steps that run after model training or periodically on live inference data.

## 30. Streamlit, Plotly Dash (Interactive Dashboarding & Web UIs for Human Oversight)

- **Role in Script Kitty:** Python frameworks for rapidly building interactive web applications and dashboards, primarily used by `Guardian` for its `Human Oversight & Feedback Interface`, and potentially by `Foundry` for internal MLOps dashboards.
- **Microscopic Details:**
  - **Rapid Prototyping:** Allow for quick development of interactive UIs directly from Python code, making it easy for `Guardian` to expose its insights to human operators.
  - **Interactive Controls:** Support widgets like sliders, buttons, dropdowns, and text inputs, enabling human operators to filter data, trigger actions (e.g., approve a policy, provide feedback), or inspect specific events.
  - **Data Visualization:** Integrate seamlessly with `Matplotlib`, `Seaborn`, and `Plotly` to display real-time metrics, logs, traces, model evaluation results, and bias detection reports.
  - **Use Cases in Script Kitty (Guardian):**
    - **Policy Review Dashboard:** Presents proposed `OPA` policy changes or violations for human review and approval/override.
    - **Performance Monitoring UI:** Displays aggregated metrics, error logs, and alerts in an actionable format.
    - **Bias Reporting Dashboard:** Visualizes detected biases in models or data.
    - **Feedback Collection Interface:** Allows human operators to provide structured feedback on Script Kitty's responses, task completions, or errors, feeding directly into `Guardian`'s learning loop.
    - **Model Debugging UI:** Allows humans to input specific queries and see how different `Skills` models or `Nexus` LLMs would respond.

## 31. Rasa Open Source (Conversational AI Framework - Potential for Structured Feedback)

- **Role in Script Kitty:** While `Nexus` handles primary chat, `Rasa` could be integrated into `Guardian`'s `Human Oversight & Feedback Interface` to provide a more structured and conversational way for humans to provide feedback to Script Kitty.
- **Microscopic Details:**
  - **NLU & Dialogue Management:** Provides its own NLU engine to understand human feedback in natural language and a dialogue management system to guide the human through a structured feedback process.

- **Custom Actions:** Allows for defining custom Python actions that connect to `Guardian`'s internal systems (e.g., logging feedback to `Kafka`, updating `MLflow` experiment metadata, triggering an evaluation pipeline).
- **Channel Connectors:** Can connect to various messaging channels, allowing `Guardian` to collect feedback from the same platforms Script Kitty interacts with.
- **Use Case:** Instead of a static form, a human could chat with a `Guardian` sub-agent: "Hey, this response was inaccurate because X. It should have said Y." The Rasa agent would guide the conversation to collect structured data (`inaccuracy_type`, `corrected_output`, `reason`), which `Guardian` then uses to fine-tune `Nexus`'s LLM or `Skills`' models.

## 32. HashiCorp Vault (Secrets Management)

- **Role in Script Kitty:** The authoritative, centralized system for securely storing, managing, and accessing all sensitive credentials, API keys, and sensitive configuration data required by Script Kitty's various components.
- **Microscopic Details:**
  - **Secret Backend:** Stores secrets in various backends (e.g., encrypted on disk, in cloud KMS).
  - **Authentication Methods:** Supports diverse authentication methods (Kubernetes Service Account, AWS IAM, GitHub, LDAP). `Foundry` would configure Kubernetes Service Account authentication, allowing `Pods` to securely authenticate with `Vault` using their assigned service account token.
  - **Dynamic Secrets:** Crucial for security. `Vault` can generate on-demand, short-lived credentials for databases (PostgreSQL), cloud providers (S3, GCS), and other services. This significantly reduces the risk of static, long-lived credentials. For instance, `Armory` needing to access a specific cloud API will request a temporary, time-bound credential from `Vault`.
  - **Leasing & Revocation:** All secrets are "leased" with a TTL. After the lease expires, the secret is automatically revoked. `Foundry` ensures that Script Kitty's applications are designed to renew leases or request new secrets as needed.
  - **Audit Logging:** All interactions with `Vault` are meticulously audited, providing a comprehensive log of who accessed what secret, when, and from where, essential for `Guardian`'s compliance and security monitoring.
  - **Policies:** Granular access control policies (`HCL` or `JSON`) define which users or service accounts can access which secrets paths, and what operations (read, create, update, delete, sudo) they can perform. This prevents unauthorized access to critical credentials.
  - **Integration with Script Kitty:**
    - `Armory`: Retrieves cloud API keys, external service credentials for resource provisioning and specialized tool access.
    - `Skills`: Retrieves database credentials for specific data analysis tasks.

- **Core**: Might retrieve credentials for accessing specialized external knowledge bases.
- **Foundry**: Manages its own internal operational secrets (e.g., Kubernetes API tokens for certain operations).

## 33. Calico / Cilium (Container Network Interface & Network Policy)

- **Role in Script Kitty:** Implement the Container Network Interface (CNI) for Kubernetes, providing the network fabric and enforcing fine-grained network policies across Script Kitty's entire microservice architecture, ensuring strict segmentation and zero-trust security.
- **Microscopic Details:**
  - **CNI Implementation:** Provide the actual networking for Kubernetes pods, assigning IP addresses and routing traffic between pods and external networks.
  - **Network Policies:** Allow `Foundry` to define declarative network rules (Kubernetes `NetworkPolicy` resources) that dictate which `Pods` can communicate with which other `Pods` or external endpoints.
  - **Micro-segmentation:** Enforces a "zero-trust" model. By default, `Pods` cannot communicate unless explicitly allowed. This creates strong isolation between Script Kitty's components. For example:
    - `Nexus` pods can only talk to `Core`'s `gRPC` endpoint.
    - `Skills` sandboxed execution pods can only initiate outbound connections to `Feast`'s online store and specific whitelisted external APIs, but not directly to `PostgreSQL` or `Vault`.
    - `Guardian`'s policy engine is highly isolated and can only receive `input` from `Core`.
  - **Prevention of Lateral Movement:** If a `Skills` agent pod is compromised, network policies prevent an attacker from easily moving to other critical Script Kitty components (e.g., `Core`, `Vault`, `etcd`).
  - **DNS Policies (Cilium):** Cilium can apply policies based on DNS names, providing more flexible and secure outbound control.
  - **Identity-Aware Policies (Cilium):** Cilium uses BPF (Berkeley Packet Filter) to provide identity-aware network policies, meaning policies can be based on the Kubernetes `ServiceAccount` or `Pod` labels, rather than just IP addresses, making them more robust to IP changes.
  - **Observability (Cilium):** Provides deep visibility into network flows within the cluster, crucial for `Foundry`'s security team to monitor traffic patterns and detect anomalies.

## 34. Istio / Linkerd (Service Mesh)

- **Role in Script Kitty:** Service mesh solutions provide a programmable infrastructure layer for inter-service communication, adding crucial features for traffic management,

observability, and, critically, enhanced security across Script Kitty's distributed microservices.

- **Microscopic Details:**
  - **Sidecar Proxy:** Each Script Kitty service `Pod` is injected with a sidecar proxy (Envoy for Istio, Linkerd's proxy). All inbound and outbound network traffic for that service goes through this proxy.
  - **Mutual TLS (mTLS):** Automatically encrypts and authenticates all communication between Script Kitty's services (e.g., `Core` to `Armory`, `Nexus` to `Core`). This establishes a strong identity for each service and ensures that only trusted services can communicate, protecting against eavesdropping and impersonation even within the cluster.
  - **Fine-grained Access Control:** Beyond basic network policies, service meshes allow `Foundry` to define highly granular `AuthorizationPolicies` based on service identity, HTTP methods, paths, headers, and even custom attributes. For example, `Core` can be authorized to `POST` to `Skills`/`execute` endpoint, but `Nexus` cannot.
  - **Traffic Management:** Enables advanced routing capabilities:
    - **Canary Deployments:** Can precisely route a small percentage of traffic to a new version of a `Skills` agent, allowing `Guardian` to monitor its performance with real traffic before a full rollout.
    - **A/B Testing:** Direct specific users or traffic segments to different versions of a `Nexus` response synthesis model.
    - **Retries, Timeouts, Circuit Breakers:** Improves the resilience of inter-service calls, ensuring Script Kitty's overall stability in the face of transient failures.
  - **Observability:** Automatically collects metrics (latency, request volume, error rates), logs, and traces for all inter-service communication, enriching `Prometheus`, `Loki`, and `Jaeger` data.
  - **Policy Enforcement:** Centralizes policy enforcement at the proxy layer, ensuring consistency.

## 35. Falco (Container Runtime Security)

- **Role in Script Kitty:** A behavioral activity monitor designed to detect anomalous activity and potential threats at the kernel level within Script Kitty's running containers, providing real-time threat detection.
- **Microscopic Details:**
  - **System Call Monitoring:** `Falco` operates by leveraging `eBPF` (extended Berkeley Packet Filter) or kernel modules to intercept and monitor system calls (e.g., `execve`, `open`, `connect`) made by processes within Script Kitty's containers. This provides deep visibility into container runtime behavior.

- **Rule Engine:** Uses a powerful and flexible rules language (`YAML`) to define suspicious behaviors. `Foundry` configures `Falco` with rules specifically tailored for Script Kitty's environment.
- **Example Rules for Script Kitty:**
  - "Alert if a `Skills` sandboxed execution container attempts to read/write to `/etc/passwd` or `/var/run/docker.sock`." (Sandbox escape attempt)
  - "Alert if an `Armory` web scraping container spawns an unexpected shell process." (Compromise detection)
  - "Alert if `Core`'s process attempts to make outbound network connections to non-whitelisted IP ranges." (Exfiltration attempt)
  - "Alert if a `Foundry` MLOps service tries to delete critical `Kubernetes` resources."
- **Real-time Alerts:** Generates alerts (e.g., to `Guardian`'s `Policy Enforcement Engine`, to `Slack`, or to a `SIEM` system) when a rule is violated, enabling immediate response to potential security incidents.
- **Policy Enforcement Point:** While `OPA` is for pre-execution checks, `Falco` is for post-execution (runtime) threat detection, forming a crucial layer of defense.

## 36. Trivy / Clair (Container Vulnerability Scanners)

- **Role in Script Kitty:** Tools for scanning Docker images for known vulnerabilities and misconfigurations, integrated into `Foundry`'s CI/CD pipelines to prevent vulnerable components from being deployed to Script Kitty's production environment.
- **Microscopic Details:**
  - **Layer-by-Layer Scanning:** Scans each layer of a Docker image (base OS, installed packages, application dependencies) against publicly available vulnerability databases (e.g., NVD, OSV, language-specific vulnerability feeds like Python Package Index).
  - **CVE Detection:** Identifies Common Vulnerabilities and Exposures (CVEs) present in the image's components.
  - **Severity Categorization:** Categorizes vulnerabilities by severity (Critical, High, Medium, Low), allowing `Foundry` to enforce policies (e.g., "fail build if any Critical/High CVEs are detected").
  - **Dependency Scanning:** Scans for vulnerabilities in language-specific dependencies (e.g., `pip` packages for Python services, `npm` packages for Node.js services).
  - **Misconfiguration Detection:** Can detect common security misconfigurations in Dockerfiles (e.g., running as root, exposed sensitive ports).
  - **SBOM (Software Bill of Materials) Generation:** Can generate a list of all components and their versions within an image, providing transparency and aiding in supply chain security.

- **Integration with CI/CD:** Integrated as a mandatory step in `Foundry`'s `GitLab CI/CD` or `GitHub Actions` pipelines. An image with detected vulnerabilities will block the pipeline from proceeding to deployment.
- **Use Cases:** Scans all `Nexus`, `Core`, `Armory`, `Skills`, `Guardian`, and `Foundry` service images before they are pushed to the `Harbor` container registry or deployed to `Kubernetes`.

## 37. Keycloak (Identity & Access Management)

- **Role in Script Kitty:** An open-source Identity and Access Management (IAM) solution providing centralized authentication and authorization services for human users interacting with Script Kitty's administrative interfaces and for internal service-to-service authentication where traditional `Kubernetes RBAC` is insufficient or requires external identity.
- **Microscopic Details:**
  - **OAuth 2.0 / OpenID Connect:** Implements industry-standard protocols for secure authentication and authorization, making it compatible with various client applications.
  - **User & Role Management:** Manages users, groups, and roles. `Foundry` defines roles for different types of human operators (e.g., "Script Kitty Admin", "Guardian Policy Reviewer", "MLOps Engineer") and assigns them appropriate permissions.
  - **Single Sign-On (SSO):** Enables human users to log in once to `Keycloak` and then access multiple Script Kitty administrative applications (e.g., `Grafana`, `MLflow UI`, `Jaeger UI`, `Guardian`'s feedback dashboard) without re-authenticating.
  - **Service Account Integration:** `Keycloak` can manage service accounts for internal Script Kitty services that need to interact with `Keycloak`-protected resources or external systems requiring OAuth tokens.
  - **Integration with Kubernetes:** Can integrate with Kubernetes for advanced authentication scenarios (e.g., authenticating users to the Kubernetes API server itself).
  - **Multi-Factor Authentication (MFA):** Supports various MFA methods, adding an extra layer of security for administrative access.
  - **Audit Logging:** Logs all authentication and authorization events, providing a clear audit trail for security compliance.
  - **Centralized Authorization:** Can serve as a Policy Enforcement Point (PEP) for authorization decisions for internal API calls that require complex external identity verification beyond `Kubernetes RBAC`.

---

**V. Data Types & Formats**

**1. Structured JSON (for Config, General Data Exchange)**

- **Role in Script Kitty:** A human-readable, widely adopted data interchange format used for configuration files, API payloads (especially where `Protobuf` is not strictly necessary or for external integrations), and storing semi-structured data.
- **Microscopic Details:**
  - **Human Readability:** Easy for developers to inspect and debug.
  - **Schema Flexibility:** More flexible than `Protobuf` for data that may evolve frequently or have optional fields, although `Pydantic` models are used in Python to enforce schema even with JSON.
  - **Use Cases:**
    - **API Payloads (REST):** For `Nexus`'s web UI API, or `Armory`'s external API calls.
    - **Configuration Files:** `ConfigMaps` in `Kubernetes` often store configurations as JSON or YAML (a superset of JSON).
    - **LLM Output:** The `Nexus LLM` and `Core LLM` are prompted to output `JSON` for structured intent and entity extraction, or plan decomposition. This structured output is validated with `Pydantic` models.
    - **Logging:** Many log messages are structured JSON for easier parsing and querying in `Loki` or `Elasticsearch`.

**2. YAML (for Kubernetes Manifests, Configuration)**

- **Role in Script Kitty:** The declarative standard for defining Kubernetes resources and other human-readable configuration files across `Foundry` and other Script Kitty components.
- **Microscopic Details:**
  - **Declarative Configuration:** `YAML` files define the *desired state* of `Kubernetes` objects (Pods, Deployments, Services, Ingresses, `CRDs`, `NetworkPolicies`, `Argo Workflows`, `Kubeflow Pipelines`). `Foundry` leverages `GitOps` to synchronize these `YAML` files with the live cluster.
  - **Readability:** Indentation-based syntax promotes readability for complex nested configurations.
  - **Use Cases:**
    - **Kubernetes Manifests:** All deployments of Script Kitty's microservices are defined as `YAML`.
    - **CI/CD Pipeline Definitions:** `GitLab CI/CD` and `GitHub Actions` pipelines are defined in `YAML`.
    - **Argo Workflows / Kubeflow Pipelines Definitions:** The DAGs and steps for ML pipelines are defined in `YAML`.
    - **DVC Pipelines (`dvc.yaml`):** Data versioning and ML pipelines are defined in `YAML`.

- **OPA Rego Policies (Data):** While Rego is the language, the input data for OPA policy evaluation is often `JSON` or `YAML`.
- **Tool Manifests (Armory):** Can be represented in `YAML` for human authoring before conversion to `Protobuf` for internal use.

## 3. HTML (Web Content)

- **Role in Script Kitty:** The raw and parsed format of web pages, primarily consumed by `Armory` for information acquisition through web scraping.
- **Microscopic Details:**
  - **Raw Source:** `Armory`'s `Playwright/Selenium` retrieve the raw HTML source of web pages.
  - **Parsed Tree:** `BeautifulSoup4`/`lxml` parse this raw `HTML` into a navigable tree structure, allowing `Armory` to precisely locate and extract data using CSS selectors or XPath.
  - **Content Extraction:** `Trafilatura`/`boilerpy3` process this `HTML` to extract clean, main content.
  - **Use Cases:** `Armory` digests this `HTML` to extract entities, facts, documents, and tool information to augment the `Global Knowledge Graph` and support `Core`'s research tasks.

## 4. Time-Series Data (Metrics)

- **Role in Script Kitty:** Numerical data points indexed by timestamps, forming the backbone of Script Kitty's performance monitoring and operational intelligence.
- **Microscopic Details:**
  - **Source:** Generated by every Script Kitty microservice and collected by `Prometheus`.
  - **Format:** Typically exposed as plain text in `Prometheus` exposition format (e.g., `my_metric_total{label="value"} 123.45`).
  - **Types:** Gauges (e.g., `active_sessions`), Counters (e.g., `requests_total`), Histograms (e.g., `request_latency_seconds_bucket`), Summaries.
  - **Storage:** Stored in `Prometheus`'s time-series database.
  - **Querying:** Analyzed using `PromQL` in `Grafana` for real-time dashboards and alerting.
  - **Use Cases:** Monitoring latency, throughput, error rates, resource utilization, LLM inference times, task success rates, and overall system health for all Script Kitty components. `Guardian` uses these for anomaly detection and triggering retraining.

## 5. Logs (Event Data)

- **Role in Script Kitty:** Structured and unstructured text messages generated by every component of Script Kitty, capturing discrete events, errors, warnings, and informational messages, crucial for debugging, auditing, and post-mortem analysis.
- **Microscopic Details:**
  - **Source:** `stdout`/`stderr` of every containerized Script Kitty service.
  - **Format:** Can be unstructured text or structured `JSON` (preferred for querying).
  - **Collection:** Collected by `Fluent Bit`/`Fluentd` agents on each `Kubernetes` node.
  - **Enrichment:** Enriched with `Kubernetes` metadata (pod name, namespace) and Script Kitty-specific `session_id`, `task_id`, `agent_name` for context.
  - **Storage:** Aggregated and stored in `Loki` (for rapid search of semi-structured logs) or `Elasticsearch` (for full-text search on structured logs).
  - **Querying:** Searched using `LogQL` in `Grafana` or `Elasticsearch` query language.
  - **Use Cases:** Troubleshooting errors, auditing actions performed by `Core` or `Skills`, tracing user interactions, and providing insights into `Guardian`'s policy decisions.

## 6. Traces (Distributed Request Flow Data)

- **Role in Script Kitty:** A sequence of spans that represent the end-to-end execution path of a single request or task across multiple microservices, providing deep visibility into distributed system behavior.
- **Microscopic Details:**
  - **Source:** Generated by `OpenTelemetry SDKs` instrumented in each Script Kitty service.
  - **Components:** Consist of `spans` (representing individual operations) linked by `trace_ids` and `parent_span_ids`.
  - **Metadata:** Each `span` contains `start/end time`, `duration`, `attributes` (e.g., `http.method`, `db.statement`, `agent.name`), and optionally `logs`.
  - **Propagation:** `Trace context` (trace_id, span_id) is propagated across network boundaries (HTTP headers, gRPC metadata) to link distributed operations.
  - **Storage:** Collected by `OpenTelemetry Collectors` and stored in `Jaeger`/`Zipkin`.
  - **Visualization:** Viewed as flame graphs or Gantt charts in `Jaeger`/`Zipkin` UI.
  - **Use Cases:** Diagnosing latency issues in `Core`'s planning, understanding complex multi-agent interactions (e.g., how `Core` orchestrates `Armory`'s research and `Skills`' code generation for a user query), and debugging distributed failures.

**7. Vector Embeddings (Numerical Representations of Semantics)**

- **Role in Script Kitty:** High-dimensional numerical vectors that capture the semantic meaning of text, code, or other data, enabling efficient semantic search and retrieval (RAG) within Script Kitty's knowledge and tool systems.
- **Microscopic Details:**
  - **Generation:** Generated by specialized `embedding models` (e.g., `sentence-transformers` models, or commercial embedding APIs) from raw text, code, or structured descriptions.
  - **High Dimensionality:** Typically hundreds or thousands of dimensions.
  - **Semantic Proximity:** Vectors that are semantically similar are closer to each other in the vector space.
  - **Storage:** Stored in `Weaviate`/`Chroma` vector databases, along with associated `metadata` and the original `text chunk`.
  - **Use Cases:**
    - **Core:** Semantic search over `GKGS` nodes for relevant knowledge.
    - **Armory:** Semantic search over `Tool Manifests` (tool descriptions) and `digested web content` (research papers, API docs) for relevant tools or information.
    - **Skills:** Semantic search for similar code examples or problem-solving approaches.
    - **RAG (Retrieval Augmented Generation):** Critically, `Nexus`, `Core`, and `Skills` use these embeddings to retrieve relevant context from the vector database, which is then fed into their `LLMs` to enhance factual accuracy and reduce hallucination.

**8. PDDL-like Syntax / Custom Rule Syntax (for Symbolic Planning & Policies)**

- **Role in Script Kitty:** A formal, declarative syntax for defining planning domains, goals, and actions for `Core`'s symbolic `HTN Planning Engine` fallback, and for `Guardian`'s custom rule sets.
- **Microscopic Details:**
  - **PDDL (Planning Domain Definition Language):** A standard language for specifying AI planning problems. While `Core`'s primary planner is LLM-driven, a simplified `PDDL-like` syntax allows for:
    - **Formal Verification:** Plans can be formally verified against the domain model.
    - **Deterministic Behavior:** Provides a predictable execution path for critical sub-tasks.
    - **Knowledge Representation:** Defines `predicates`, `actions` (with `parameters`, `preconditions`, `effects`), and `goals`.

- - **Custom Rule Syntax:** A custom, simplified declarative language for defining specific operational rules or heuristics that `Core` might use, or for `Guardian`'s lightweight ethical rules.
  - **Parsing:** `ply` (Python Lex-Yacc) is used to parse these syntaxes into an Abstract Syntax Tree (AST) that `Core` or `Guardian` can interpret.
  - **Use Cases:**
    - **Core:** For highly critical or well-defined planning sub-problems, the `LLM` might generate a `PDDL-like` plan which is then parsed and verified by the symbolic planner.
    - **Guardian:** For defining specific, non-negotiable safety rules or operational guardrails that are too critical or simple to be managed by a full `OPA` policy or an `Ethical LLM`.
  - **Integration with GKGS:** The knowledge encoded in these symbolic rules/plans can be stored as structured data within `Core`'s `Global Knowledge Graph`.

## 9. Python Code / Shell Scripts (for Agent Execution)

- **Role in Script Kitty:** The executable artifacts generated by `Skills` agents or retrieved by `Armory` for execution within sandboxed environments.
- **Microscopic Details:**
  - **Skills Agent Generation:** `Skills` agents (`Code-focused LLMs`) generate Python code for data analysis, mathematical computations, or specific task automation.
  - **Armory Tool Wrappers:** `Armory` can generate Python scripts or shell scripts as wrappers for external CLI tools or APIs.
  - **Sandboxed Execution:** All such code is executed within isolated, ephemeral Docker containers (managed by `Foundry`) with strict resource limits and network policies (enforced by `Calico`/`Cilium`). `Skills` agents use `Jupyter Kernel Gateway` or custom `subprocess` calls within these containers.
  - **Security Context:** Containers run with minimal privileges (`non-root user`, `seccomp-bpf` filters, `AppArmor`/`SELinux` profiles) to mitigate risks of code execution.
  - **Input/Output:** Code takes structured inputs (from `Core`) and produces structured outputs (e.g., `JSON`, `CSV`, images) back to `Core`.
  - **Testing:** `Skills` can generate unit tests (`pytest`) for its own generated code and execute them within the sandbox to verify correctness before returning results.

## 10. IaC Templates (Terraform HCL, Ansible Playbooks)

- **Role in Script Kitty:** Declarative Infrastructure as Code templates used by `Armory` for provisioning computational resources (VMs, specialized GPUs, storage buckets) in cloud environments.
- **Microscopic Details:**
  - **Terraform HCL (HashiCorp Configuration Language):**
    - **Declarative:** Defines the desired state of infrastructure (e.g., "I need a VM with 4 CPUs, 16GB RAM, and a specific GPU type in region X").
    - **Cloud Agnostic:** Supports various cloud providers (AWS, GCP, Azure) and on-premise infrastructure.
    - **State Management:** `Terraform` maintains a state file (often stored in `S3`/`MinIO`) that maps the declared configuration to real-world resources.
    - **Use Cases:** `Armory` uses `Terraform` to spin up ephemeral cloud VMs for large-scale data processing or for accessing very specialized hardware not directly available in `Kubernetes`.
  - **Ansible Playbooks:**
    - **Configuration Management:** Used for configuring software, installing dependencies, and managing services on provisioned VMs.
    - **Idempotent:** Running a playbook multiple times has the same effect, ensuring consistent configuration.
    - **Agentless:** Communicates over SSH, requiring no agent on the target machine.
    - **Use Cases:** Once `Terraform` provisions a VM, `Armory` can use `Ansible` to install specific libraries, drivers, or software versions required for a specialized tool or quantum SDK.
  - **Secure Parameters:** `Armory` dynamically injects sensitive parameters (e.g., API keys, resource IDs) from `HashiCorp Vault` into these templates at runtime, rather than hardcoding them.
  - **Auditing:** All `IaC` deployments are logged and auditable, contributing to `Guardian`'s oversight.

---

**VI. Expert Addition: Critical Missing Resources/Tools**

Based on the highly detailed and ambitious scope of Script Kitty, a crucial aspect not explicitly detailed but implicitly required for its practical operation and continuous evolution is a **Centralized Source Code Repository and Version Control System** beyond just Git's core capabilities for configuration.

**1. Git-based Source Code Management (e.g., GitLab, GitHub, Gitea - as a Platform)**

- **Reason for Inclusion:** While "Git commits" are mentioned in `Foundry`'s CI/CD section and `DVC`'s integration, the underlying platform that *hosts* these Git repositories and

provides collaborative features, code review workflows, and integrated CI/CD runners is fundamental for building, managing, and evolving an AGI of this complexity. Without a robust, centralized Git platform, the collaborative development, continuous integration, and declarative deployments described would be impractical.

- **Microscopic Details:**
  - **Centralized Repository Hosting:** Provides a secure, version-controlled home for all Script Kitty's source code: `Nexus` microservices, `Core`'s planning logic, `Armory`'s scraper code, `Skills`' model code, `Guardian`'s policy definitions and evaluation scripts, and `Foundry`'s infrastructure code (`IaC`, `Kubernetes` manifests).
  - **Version Control (Git):** At its core, it leverages Git for distributed version control, enabling branching, merging, pull requests, and commit history tracking. Every change to Script Kitty's intelligence or infrastructure is auditable.
  - **Collaborative Development:** Supports features like:
    - **Pull Requests / Merge Requests:** Essential for code review processes, ensuring code quality, security, and alignment with architectural principles before merging into `main` branches. `Foundry`'s CI/CD pipelines would be triggered by these.
    - **Issue Tracking:** Integrated issue trackers allow for managing bugs, features, and tasks related to Script Kitty's development.
    - **Wiki/Documentation:** Provides centralized documentation for Script Kitty's architecture, API specifications, development guidelines, and operational procedures.
  - **Integrated CI/CD Runners:** Platforms like `GitLab CI/CD` and `GitHub Actions` provide built-in or self-hosted runners that execute `Foundry`'s defined CI/CD pipelines directly in response to Git events, providing the compute for building, testing, and deploying Script Kitty.
  - **Access Control:** Robust authentication (often integrated with `Keycloak`) and authorization (per-repository, per-branch permissions) to control who can read, write, or merge code.
  - **Webhooks:** Triggers `Foundry`'s `CI/CD` pipelines or other automation in response to specific Git events (e.g., push, pull request, tag creation).
  - **Code Search & Navigation:** Facilitates navigating the vast codebase of Script Kitty.
  - **Security Features:** Often includes features like static application security testing (SAST), dependency scanning (complementing `Trivy`/`Clair`), and secret detection within repositories, adding another layer of defense.
  - **Use Cases across Script Kitty:**
    - **Foundry:** The Git repository is the single source of truth for all `IaC` and application deployments. `Argo CD`/`Flux CD` continuously sync from these repos.

- **Guardian:** Policy definitions (Rego), evaluation scripts, and bias detection logic are stored and versioned here.
- **Core, Nexus, Armory, Skills:** All their application code, LLM fine-tuning scripts, data processing logic, and model definitions reside here.
- **DVC:** `.dvc` files that point to the actual data in object storage are committed to these Git repositories.

# Workflow

**Phase 1: User Input Ingestion and Initial Understanding (Script Kitty.Nexus)**

This phase initiates the journey, transforming raw user intent into a structured, actionable request.

1. **User Input Reception (Via ChannelAdapter):**

   ○ **Minute Detail:** The user types their request into a chosen interface (e.g., a web chat UI, Slack, Telegram). The specific `ChannelAdapter` (e.g., `WebChatAdapter` via WebSocket, `SlackAdapter` via Events API webhook) configured for that channel receives the raw text: `"Hey, can you research the latest advancements in quantum error correction, summarize key breakthroughs, and then propose a Python code snippet that simulates a basic Shor's algorithm for a 3-qubit system?"`

   ○ **Technical Dissection:**
      ■ **Protocol:** HTTPS POST (for webhooks) or WebSocket frame (for direct web chat).
      ■ **ChannelAdapter Function:** `receive_raw_input(channel_id, user_id, raw_text)`.
      ■ **Data Transformation:** The `ChannelAdapter` immediately wraps the raw input into a `ScriptKittyMessageSchema` Protobuf message.
         ■ `script_kitty_message.user_id = <unique_user_id>`
         ■ `script_kitty_message.session_id = <existing_or_new_session_id>` (retrieved from a lightweight local cache or by querying Session Manager if new)
         ■ `script_kitty_message.timestamp = <current_utc_timestamp>`
         ■ `script_kitty_message.channel_type = 'WEB_CHAT'` (or 'SLACK', 'TELEGRAM')
         ■ `script_kitty_message.raw_text = "..."` (the full user query)
         ■ `script_kitty_message.event_type = 'USER_INPUT'`
         ■ `script_kitty_message.message_id = <unique_message_id>`
      ■ **Output:** The fully structured `ScriptKittyMessageSchema` is then queued internally within Nexus for the next processing step or directly passed via a local function call.

2. **Input Normalization & Security Filtering:**

- **Minute Detail:** The `ScriptKittyMessageSchema` (specifically its `raw_text` field) passes through a series of deterministic pre-processing steps.
- **Technical Dissection:**
  - **Function:** `normalize_and_filter_input(script_kitty_message)`.
  - **Steps (Sequential Application):**
    - **Unicode Normalization:** `unicodedata.normalize('NFC', script_kitty_message.raw_text)` is applied to ensure consistent character representation.
    - **Whitespace & Punctuation Cleanup:** Regular expressions (e.g., `re.sub(r'\s+', ' ', text)`) collapse multiple spaces, remove leading/trailing whitespace.
    - **Lightweight Spell Correction:** A fast, approximate string matching algorithm (e.g., `pyspellchecker` or custom `SymSpell`-based module) corrects obvious typos (e.g., "qunatum" to "quantum").
    - **Heuristic PII Redaction:** Basic regex patterns (e.g., for email addresses `\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b`) identify and redact common PII, replacing it with `[REDACTED_PII]`. (Crucially, this is a preliminary step; Guardian performs a more robust, policy-driven redaction).
    - **Heuristic Toxicity Scoring:** A small, pre-trained `scikit-learn` classifier (e.g., `LogisticRegression` trained on a corpus of toxic/non-toxic phrases) or a simple keyword matching system rapidly assigns a preliminary toxicity score. If a high score, a flag (`is_flagged_toxic: true`) is added to the `script_kitty_message`.
    - **Anti-Spam/Flood Control:** A rate limiter (e.g., using Redis counters per `user_id` and `channel_type`) checks for excessive requests within a time window. If triggered, the message is dropped, and an internal alert is logged.
  - **Output:** An updated `script_kitty_message` with cleaned `raw_text` and relevant flags.

3. **Session Context Retrieval (Session Manager):**

   - **Minute Detail:** Before NLU, Nexus retrieves the ongoing conversational context for the `session_id` to provide the LLM with relevant history.
   - **Technical Dissection:**
     - **gRPC Call:** Nexus (as a gRPC client) initiates a call to the `SessionManager` service.

- **Endpoint:**
  `SessionManager.GetSessionContext(session_id=<script_kitty_message.session_id>)`.
- **SessionManager Logic:**
  - **Cache Check (Redis):** Attempts to fetch session data from Redis (`redis_client.get(session_id)`). If found and `last_activity_timestamp` is recent, it returns the cached `session_context_json`.
  - **DB Fallback (PostgreSQL):** If not in Redis or if expired, it queries PostgreSQL (`psycopg2` or ORM like SQLAlchemy): `SELECT context FROM sessions WHERE session_id = <session_id>`.
  - **New Session Creation:** If no session exists, `SessionManager` creates a new entry in both Redis and PostgreSQL with an empty context and a new `session_id`.
- **Data Structure:** The `session_context` is a JSON object with fields like `conversation_history` (list of past `ScriptKittyMessageSchema` objects or summarized turns), `extracted_entities`, `user_preferences`, `current_task_state`.
- **Output:** The `session_context_json` is returned to Nexus.

4. **Core NLU & Intent Routing (Nexus LLM & Router):**

- ○ **Minute Detail:** The core understanding phase. The Nexus LLM processes the user query and session context to determine intent and extract entities, guiding the routing decision.
- ○ **Technical Dissection:**
  - **Function:** `analyze_and_route(script_kitty_message, session_context)`.
  - **LLM Inference Request (vLLM/Ollama):** Nexus's LLM Inference Client (using `grpc.aio` or HTTP client) sends a request to the `NexusLLMService` (running `vLLM` or `Ollama` for inference).
    - **Request Payload:**
      - **System Prompt:** "You are an expert NLU engine for Script Kitty. Extract the user's primary intent and all relevant entities. Output a JSON object. Available intents: `['RESEARCH_AND_SUMMARIZE', 'CODE_GENERATION', 'GENERAL_CHAT', 'TASK_MANAGEMENT']`. Entities: `['topic', 'algorithm', 'system_size', 'language']`. Ensure all fields are present, even if empty. If unsure, default to `GENERAL_CHAT`."

- **Context:** `session_context` (e.g., `{"conversation_history": "...", "user_preferences": "..."}`).
- **User Input:** `script_kitty_message.raw_text`.
- **Inference Parameters:** `temperature=0.2` (for deterministic output), `max_tokens=150`, `stop_sequences=['\n\n']`.

**Expected LLM Output (JSON):**
JSON

```json
{
  "intent": "RESEARCH_AND_SUMMARIZE_AND_CODE_GENERATION",
  "entities": {
    "research_topic": "latest advancements in quantum error correction",
    "summarize_aspects": "key breakthroughs",
    "code_task": "simulate a basic Shor's algorithm",
    "algorithm": "Shor's algorithm",
    "system_size": "3-qubit system",
    "language": "Python"
  },
  "confidence": 0.95
}
```

- ■
- **JSON Schema Validation:** The LLM's raw JSON output is immediately validated against a `ScriptKittyIntentSchema` (Protobuf-defined). If validation fails, a fallback prompt might be tried or a `ParsingError` logged.
- **Semantic Router Logic:**
  - **Decision Tree/Rule-based on Intent:** Based on `parsed_intent.intent`:
    - If `RESEARCH_AND_SUMMARIZE_AND_CODE_GENERATION` (complex, multi-modal), the request is *always* escalated to **Script Kitty.Core**.
    - If `GENERAL_CHAT` and `confidence > 0.8`, handled by Nexus's own response generation.
    - **Fallback (Pre-trained Classifier):** If LLM confidence is low or a very common, simple intent (e.g., "hello") is detected, a fine-tuned `DistilBERT` model (loaded from Hugging Face Transformers) might cross-check the intent for faster, higher-confidence routing.

- **Enrichment:** The `script_kitty_message` is enriched with `script_kitty_message.parsed_intent` and `script_kitty_message.extracted_entities`.
- **Output:** A decision to either:
    - Proceed to Nexus's Response Synthesis (for simple chat).
    - **Escalate to Script Kitty.Core** (for complex tasks like this scenario).

5. **Escalation to Script Kitty.Core (gRPC Communication):**

- **Minute Detail:** The fully enriched `ScriptKittyMessageSchema` is now handed over to the central orchestrator.
- **Technical Dissection:**
    - **gRPC Call (Asynchronous):** Nexus (acting as a gRPC client) initiates an asynchronous remote procedure call to **Script Kitty.Core**.
        - **Endpoint:** `CoreService.ProcessUserRequest(script_kitty_message)`.
    - **Payload:** The `script_kitty_message` Protobuf, now containing the raw text, session info, and the `parsed_intent` with `extracted_entities`.
    - **Error Handling:** Nexus implements retry logic with exponential backoff for gRPC connection errors or Core service unavailability, ensuring messages are not lost. A timeout is also applied; if Core doesn't acknowledge within a specific time, Nexus logs an internal error and might send a "System Busy" message to the user.
    - **Output:** The request is successfully passed to Core, and Nexus awaits a response via gRPC stream or subsequent updates. Nexus updates the `session_context` (via `SessionManager`) to indicate that a task has been sent to Core and is pending.

---

**Phase 2: Strategic Planning and Orchestration (Script Kitty.Core)**

Core receives the structured request and breaks it down into a hierarchical, executable plan, dynamically orchestrating specialized agents.

1. **Request Ingestion & Validation (Core's Internal Communication Hub):**

- **Minute Detail:** Core's gRPC server receives the `ProcessUserRequest` call from Nexus.
- **Technical Dissection:**

- **gRPC Server Listener:** Core runs a gRPC server that listens for incoming Nexus requests.
- **Input Validation:** The incoming `script_kitty_message` Protobuf is validated against its schema to ensure integrity. Malformed requests are rejected.
- **Internal Queueing:** The request is immediately placed into an internal, high-priority message queue (e.g., a Kafka topic or an in-memory buffered channel) for asynchronous processing, preventing the gRPC endpoint from blocking.
- **Acknowledgment:** Core sends an immediate gRPC acknowledgment back to Nexus (a simple `Success` status or `Processing` status) indicating receipt, allowing Nexus to update its internal state.

2. **Goal Decomposition & HTN Planning (Core's Planning Engine):**

- **Minute Detail:** This is the core's strategic brain at work, transforming a high-level user goal into a detailed, executable plan.
- **Technical Dissection:**
  - **LLM for High-Level Planning (e.g., Llama 3-70B, or GPT-4o via API):** A more powerful LLM than Nexus's, hosted either locally via vLLM on a high-end GPU cluster (Foundry-managed) or accessed via a secure, rate-limited API gateway (Foundry-managed).
    - **System Prompt:** "You are Script Kitty's expert planning engine. Decompose the user's request into a Hierarchical Task Network (HTN) of sub-goals and atomic actions. Each action must specify the `agent_type` (Armory, Skills, Guardian), `function_name`, and `parameters`. Output a JSON array of sequential steps, with dependencies. Consult the available tool manifests from the Global Knowledge Graph. Break down the task 'research quantum error correction, summarize, then simulate Shor's algorithm'."
    - **Context:** `script_kitty_message.parsed_intent`, `script_kitty_message.extracted_entities`, relevant `session_context`, and crucially, dynamically retrieved tool manifests from the **Global Knowledge Graph (GKGS)**. The GKGS provides information on what tools (from Armory) and skills (from Skills) are available and their capabilities.
    - **Planning Algorithm:** The LLM generates a proposed `PlanSchema` (Protobuf-defined JSON array of steps).

**Generated Plan (Example `PlanSchema`):**
 JSON
[

```
 {"step_id": "s1", "agent_type": "Armory", "function_name": "research_web_and_digest",
"parameters": {"query": "latest advancements in quantum error correction", "keywords":
["breakthroughs", "challenges"], "output_format": "structured_documents"}, "dependencies": []},
  {"step_id": "s2", "agent_type": "Skills", "function_name": "summarize_documents",
"parameters": {"documents_source": "s1.output.structured_documents", "length": "concise",
"focus": "key_breakthroughs"}, "dependencies": ["s1"]},
  {"step_id": "s3", "agent_type": "Skills", "function_name": "generate_code_snippet",
"parameters": {"task_description": "simulate basic Shor's algorithm", "algorithm": "Shor's
algorithm", "qubit_count": 3, "language": "Python", "context_summary":
"s2.output.summary_text"}, "dependencies": ["s2"]},
  {"step_id": "s4", "agent_type": "Skills", "function_name": "execute_sandboxed_code",
"parameters": {"code_snippet": "s3.output.code_snippet", "language": "Python",
"expected_output_format": "execution_result"}, "dependencies": ["s3"]},
  {"step_id": "s5", "agent_type": "Core", "function_name": "synthesize_final_response",
"parameters": {"research_summary": "s2.output.summary_text", "simulation_result":
"s4.output.execution_result"}, "dependencies": ["s2", "s4"]}
]
```

- ■
  - **Symbolic HTN Planner (Fallback/Verification - Optional but Recommended):** For critical or highly constrained tasks, a traditional symbolic HTN planner (e.g., `pyddl` or custom implementation) can formally verify the LLM's plan against predefined domain rules or even generate plans in well-defined problem spaces. This adds a layer of robustness.
  - **GKGS Integration (RAG):** During planning, Core dynamically retrieves relevant information from the GKGS (using RAG pipelines with `LlamaIndex` or `Haystack` over Weaviate/Chroma vector store) to inform planning decisions (e.g., available tool parameters, past success rates of certain Skills agents).

3. **Plan Validation & Safety Check (Integration with Script Kitty.Guardian):**

   - ○ **Minute Detail:** Before any action is taken, the generated plan is submitted to Guardian for policy enforcement.
   - ○ **Technical Dissection:**
     - **gRPC Call:** Core initiates a gRPC call to **Script Kitty.Guardian**.
       - **Endpoint:** `GuardianService.ReviewPlan(plan_schema, session_id, user_id)`.
     - **Guardian's Policy Enforcement Engine:**
       - **OPA (Open Policy Agent):** Guardian evaluates the `plan_schema` against predefined security, ethical, and operational policies written in Rego language.

- **Ethical LLM:** A specialized LLM within Guardian might analyze the semantic meaning of the plan for potential ethical dilemmas (e.g., "is researching a company's internal network allowed?").
- **Decision:** Guardian returns a `PolicyDecision` Protobuf (`ALLOW`, `DENY`, `REQUIRES_HUMAN_INTERVENTION`).

- **Core's Reaction:**
  - If `ALLOW`: Core proceeds to execution.
  - If `DENY`: Core logs the violation, informs Nexus to send an error message to the user, and updates the `session_context` with the denial.
  - If `REQUIRES_HUMAN_INTERVENTION`: Core pauses execution, sends an alert to Guardian's human oversight interface, and informs Nexus to tell the user that the request requires review. (Guardian manages the human approval process).

4. **Agent Orchestration & Lifecycle Management:**

- **Minute Detail:** Core initiates the execution of the validated plan, activating and monitoring the specialized agents.
- **Technical Dissection:**
  - **Workflow Orchestration Engine (Internal or Kubeflow/Argo Workflows):** Core's internal state machine (or a dedicated workflow engine like Kubeflow Pipelines/Argo Workflows managed by Foundry) begins executing the `PlanSchema` step by step.
  - **Dynamic Agent Activation (Kubernetes API):** For each step in the plan, Core determines the `agent_type` and `function_name`.
    - **If agent already running:** Core sends a gRPC task request to the existing agent instance.
    - **If agent not running/requires specific resources:** Core (via Foundry's Kubernetes Python Client) interacts with the Kubernetes API to dynamically deploy or scale up the necessary agent service.
      - `kubernetes.client.AppsV1Api.create_namespaced_deployment(namespace='script-kitty-skills', body=skills_deployment_manifest)`.
      - The `skills_deployment_manifest` would include specific resource requests (e.g., GPU count), container image from Foundry's registry, and environment variables.
  - **Task Dispatch (gRPC):** Once the target agent is confirmed active, Core sends a specific gRPC task request.
    - **Example (Step 1 - Armory):**
      `ArmoryService.ExecuteResearch(research_query="...", output_format="structured_documents",`

```
                      task_id=<unique_task_id>,
                      session_id=<session_id>).
```
- **Asynchronous Handling:** Core initiates these gRPC calls asynchronously and stores a mapping of `task_id` to `session_id` and `plan_step_id` to await results.
  - **Progress Updates to Nexus:** Core sends periodic `PROGRESS_UPDATE` messages back to Nexus via its gRPC stream/callback, allowing Nexus to provide real-time updates to the user (e.g., "Core is researching quantum error correction...").

---

**Phase 3: Specialized Task Execution (Script Kitty.Armory & Script Kitty.Skills)**

The execution layer, where the heavy lifting of research, data digestion, code generation, and sandboxed execution occurs.

**A. Armory's Role (for Step 1):**

1. **Research & Digest Execution (Armory):**
   - **Minute Detail:** Armory receives the `ExecuteResearch` task from Core and initiates web scraping and data processing.
   - **Technical Dissection:**
     - **gRPC Listener:** Armory's gRPC server receives `ExecuteResearch` request.
     - **Dynamic Resource Allocation (Foundry):** If the research requires significant compute (e.g., heavy browser automation), Armory might request ephemeral Kubernetes Pods from Foundry via Foundry's Kubernetes API, ensuring isolated and scaled execution.
     - **Headless Browser Automation (Playwright/Selenium):** Armory launches a headless browser instance (`playwright.sync_api.playwright.chromium.launch()`) in a dedicated, sandboxed container (isolated by Foundry's Kubernetes network policies and potentially Firejail).
     - **Navigation & Scraping:** Navigates to search engines (Google Scholar, arXiv) or specific research portals using the `query` from Core. It executes JavaScript, interacts with DOM elements (`page.click()`, `page.fill()`), and extracts HTML content (`page.content()`).
     - **HTML Parsing (BeautifulSoup4/lxml):** The raw HTML is parsed efficiently.
     - **Text Extraction (Trafilatura/boilerpy3):** Focuses on extracting clean, relevant text content, filtering out boilerplate.

- **LLM for Contextual Extraction/Summarization:** A dedicated Armory LLM (e.g., fine-tuned Llama 3-8B) is used to process large chunks of text extracted from web pages.
    - **Prompt:** "Summarize the key breakthroughs in quantum error correction from the following document, focusing on novel techniques and experimental validations. Extract any mentioned key figures or institutions."
    - **Output:** Structured JSON containing summaries, identified entities (e.g., research papers, authors, institutions), and links to original sources.
- **Data Validation & Structuring (Pydantic):** The extracted data is validated against a predefined schema to ensure consistency.
- **Output Data Storage (Foundry's Data Lake):** The `structured_documents` are stored in Foundry's S3/MinIO data lake (e.g., `s3://armory-research-docs/quantum-error-correction-2025.json`) with appropriate versioning (DVC).
- **Progress Updates to Core:** Armory sends `TASK_UPDATE` messages to Core's gRPC endpoint with status like `RESEARCH_IN_PROGRESS`, `PROGRESS_PERCENTAGE=...`.
- **Task Completion:** Once research is complete, Armory sends a `TASK_RESULT` message to Core.
    - **Payload:** `task_id`, `status='COMPLETED'`, `output_data={'structured_documents_uri': 's3://...', 'summary_preview': '...'}`.

## B. Skills' Role (for Steps 2, 3, 4):

2. **Summarize Documents Execution (Skills - Step 2):**

   - **Minute Detail:** Skills retrieves the research output from Armory and generates a concise summary.
   - **Technical Dissection:**
        - **gRPC Listener:** Skills' gRPC server receives `ExecuteSummarize` request from Core, with `documents_source` pointing to the S3 URI from Armory's output.
        - **Data Retrieval:** Skills (via Foundry's S3 SDK client) fetches the `structured_documents` from the specified S3 URI.
        - **Specialized LLM for Summarization:** A fine-tuned summarization LLM (e.g., BART or T5) within the Skills service, optimized for condensing technical documents.

- **Input:** The retrieved `structured_documents` (concatenated or processed in chunks).
- **Prompt:** "Based on the provided research documents on quantum error correction, generate a concise summary focusing on key breakthroughs and experimental validations."
- **Output:** `summary_text`.
- **Task Completion:** Skills sends a `TASK_RESULT` to Core with `output_data={'summary_text': '...'}`.

3. **Generate Code Snippet Execution (Skills - Step 3):**

   - **Minute Detail:** Skills uses its code generation capabilities to propose a Python simulation for Shor's algorithm.
   - **Technical Dissection:**
     - **gRPC Listener:** Skills receives `ExecuteCodeGeneration` request from Core, with `task_description`, `algorithm`, `qubit_count`, `language`, and the `context_summary` from the previous step.
     - **Code-focused LLM (e.g., CodeLlama-34B, DeepSeek Coder):** A specialized generative LLM, specifically fine-tuned for code generation and understanding programming contexts.
       - **System Prompt:** "You are an expert Python programmer assisting in quantum computing simulations. Generate a basic, runnable Python code snippet to simulate Shor's algorithm for a 3-qubit system. Focus on demonstrating the core principles rather than full optimization. Use Qiskit. Include comments. Here is context: [summary_text from s2]."
       - **Output:** The LLM generates the Python code snippet as a string.
     - **Syntax/Basic Semantic Validation:** Before returning, the generated code is passed through a lightweight syntax parser (e.g., Python's `ast` module) and a basic semantic linter to catch obvious errors. If errors, the LLM might be prompted to retry.
     - **Output:** `code_snippet` string.
     - **Task Completion:** Skills sends a `TASK_RESULT` to Core with `output_data={'code_snippet': '...'}`.

4. **Execute Sandboxed Code Execution (Skills - Step 4):**

   - **Minute Detail:** The generated Python code is executed in a highly isolated, secure environment to get its output without risking the core system.
   - **Technical Dissection:**
     - **gRPC Listener:** Skills receives `ExecuteSandboxedCode` request from Core, with the `code_snippet` and `language`.
     - **Ephemeral Container Provisioning (Foundry's Kubernetes API):** Skills (via Foundry's Kubernetes client) requests a new, ephemeral

Docker container (e.g., `python-sandbox-env`) within a dedicated, restricted Kubernetes namespace (e.g., `script-kitty-sandbox`).

- **Container Configuration:**
    - **Minimal Base Image:** Built from a hardened, minimal base image (e.g., `distroless` for Python, or Alpine Linux with necessary libs).
    - **Resource Limits:** Strict CPU (`limits.cpu: 500m`), memory (`limits.memory: 256Mi`), and potentially execution time limits (`activeDeadlineSeconds` on the Kubernetes Job).
    - **Network Policies:** Zero-egress network policies (via Calico/Cilium) prevent the container from making any outbound network connections.
    - **Read-Only Root Filesystem:** Prevents writing to system directories.
    - **Security Context:** `runAsNonRoot: true`, `allowPrivilegeEscalation: false`, `seccompProfile: {type: RuntimeDefault}`.
    - **seccomp-bpf/AppArmor:** Custom seccomp profiles allow only a strict whitelist of necessary system calls, preventing malicious actions.
    - **Firejail/Bubblewrap (Optional Layer):** For extreme isolation, these user-space sandboxing tools can be run *inside* the Docker container for an extra layer of defense against container escapes.
- **Code Injection & Execution:**
    - The `code_snippet` is mounted as a temporary read-only file within the sandbox container.
    - The `subprocess` module (`subprocess.run(["python", "/tmp/script.py"], capture_output=True, timeout=execution_timeout)`) is used to execute the script.
    - **Output Capture:** `stdout` and `stderr` are captured.
- **Result Transmission:** The captured output (`execution_result`, `error_output`) is transmitted back to Skills.
- **Sandbox Teardown:** The ephemeral container is immediately terminated and garbage collected by Kubernetes once execution is complete or timed out.
- **Task Completion:** Skills sends a `TASK_RESULT` to Core with `output_data={'execution_result': '...', 'errors': '...'}`. Guardian would monitor for any security policy violations detected by Falco during this step.

**Phase 4: Final Synthesis and Response Generation (Script Kitty.Core & Script Kitty.Nexus)**

Core consolidates all results, and Nexus crafts the final, persona-consistent response for the user.

1. **Synthesize Final Response (Core - Step 5):**

   ○ **Minute Detail:** Core receives all task results from Armory and Skills and combines them into a coherent final output.
   ○ **Technical Dissection:**
      ■ **Dependency Resolution:** Core's workflow engine waits for `TASK_RESULT` messages from `s2` (summarization) and `s4` (code execution) to be `COMPLETED`.
      ■ **Data Aggregation:** Core retrieves `s2.output.summary_text` and `s4.output.execution_result` (and `errors` if any).
      ■ **Core's LLM for Synthesis:** Core's planning LLM (or a smaller, dedicated synthesis LLM) can be used to weave together disparate pieces of information into a single, cohesive narrative.
         ■ **Prompt:** "Combine the following research summary on quantum error correction and the output from a Shor's algorithm simulation into a single, comprehensive response for the user. Format it clearly using Markdown, include the code if successful, and explain the simulation output. Ensure the tone is expert and helpful."
         ■ **Input:** `research_summary`, `simulation_result`, `original_user_query`.
         ■ **Output:** `final_raw_response_text` (Markdown formatted).
      ■ **Error Handling:** If any sub-task failed (e.g., code execution error), Core's synthesis LLM is prompted to explain the error gracefully to the user and suggest next steps.
      ■ **Internal Audit Logging:** The full `PlanSchema` and all `TASK_RESULT` messages are logged to Foundry's centralized logging system for auditability and future learning.
      ■ **Output:** The `final_raw_response_text` is prepared to be sent to Nexus.
2. **Response Dispatch to Nexus (gRPC):**

   ○ **Minute Detail:** Core sends the final response to Nexus for formatting and delivery.
   ○ **Technical Dissection:**

- **gRPC Call:** Core makes a gRPC call to Nexus.
  - **Endpoint:** `NexusService.DeliverAgentResponse(response_schema)`.
- **Payload:** `response_schema` (Protobuf-defined) containing:
  - `session_id`: To ensure Nexus targets the correct user.
  - `response_text`: The `final_raw_response_text` from Core.
  - `source_agent`: "Core" (or "Script Kitty").
  - `event_type`: "TASK_COMPLETED_RESPONSE".
  - `suggested_actions`: (Optional) e.g., `["Ask follow-up questions", "Save summary", "Refine code"]`.
- **Core's Internal State Update:** Core marks the `session_id`'s current task as `COMPLETED`.

3. **Persona Infusion & Final Rendering (Nexus):**

- **Minute Detail:** Nexus takes Core's synthesized response and applies the overarching "Script Kitty" persona, then dispatches it to the user's channel.
- **Technical Dissection:**
  - **Nexus LLM (Persona Layer):** While Core provides the raw content, Nexus's LLM can apply a final stylistic pass if necessary, ensuring the tone and phrasing consistently match the "empirical, unparalleled inventive, limitlessly creative" persona. For this, `response_text` is fed to Nexus LLM with a specific "persona refinement" prompt.
    - **Prompt:** "Refine the following text to embody the Script Kitty persona: empirical, inventive, creative, expert in Computer Science, AI, Programming. Add a polite and encouraging closing. Here's the text: [Core's response]."
    - **Output:** `final_persona_infused_text`.
  - **Group Chat Attribution:** Nexus prefixes the response with "The Nexus has completed your complex request, with contributions from Armory and Skills:" or "Script Kitty reports:".
  - **Markdown Rendering:** Ensures the Markdown is correctly interpreted and rendered for the specific channel.

4. **User Output Dispatch (ChannelAdapter):**

- **Minute Detail:** The final, polished response is sent back to the user.
- **Technical Dissection:**
  - **Function:** `send_output_to_user(session_id, final_persona_infused_text, suggested_actions)`.
  - **Channel-Specific Translation:** The `ChannelAdapter` for the specific `channel_type` translates the `final_persona_infused_text` and `suggested_actions` into the native format of the messaging platform

(e.g., Slack Blocks for rich formatting, Telegram Keyboard for buttons, HTML/CSS for web UI).
- **API Call:** The adapter makes the final API call to the respective platform (e.g., `slack_web_client.chat_postMessage()`, `telegram_bot.send_message()`, WebSocket message).
- **Output:** The user receives the comprehensive response in their chosen chat interface.

---

**Phase 5: Continuous Learning and Feedback Loops (Script Kitty.Guardian & Script Kitty.Foundry)**

The vital processes that enable Script Kitty to continuously improve, self-correct, and evolve its capabilities.

1. **Real-time Performance Monitoring & Anomaly Detection (Foundry & Guardian):**

   - **Minute Detail:** As the workflow executes, Foundry's observability stack (Prometheus, Loki, OpenTelemetry) constantly streams performance and operational data. Guardian actively consumes and analyzes this data.
   - **Technical Dissection:**
     - **Metrics Collection:** Prometheus scrapes `http_requests_total`, `api_call_latency_seconds`, `llm_inference_duration_milliseconds`, `agent_task_success_rate`, `sandbox_resource_utilization`, `policy_violation_count` from all Script Kitty components (Nexus, Core, Armory, Skills, Guardian itself, and Foundry services).
     - **Log Ingestion:** Fluent Bit/Fluentd pushes all logs (including detailed execution traces, LLM prompts/responses, task results) to Loki/Elasticsearch.
     - **Distributed Tracing:** OpenTelemetry traces record the full end-to-end journey of the user's request, linking all spans across services.
     - **Guardian's Anomaly Detection (Foundry's ML Pipeline):**
       - Guardian runs dedicated anomaly detection models (e.g., Isolation Forest, Autoencoders on time-series metrics from Prometheus, or log parsing for unusual patterns) on streams of performance metrics and logs.
       - **Triggers:** Anomalies like `skills_sandbox_escape_attempt`, `core_planning_time_exceeded`, `nexus_nlu_confidence_drop`, or `armory_web_scrape_rate_limit_exceeded` are immediately detected.

- **Alerting:** Guardian triggers alerts through Foundry's Grafana/Alertmanager, notifying human operators and internal Script Kitty components (e.g., Core might automatically reduce load if a bottleneck is detected).

2. **Automated Evaluation & Training Trigger (Guardian & Foundry):**

- ○ **Minute Detail:** Based on successful task completions, failures, or detected anomalies, Guardian evaluates performance and determines if model retraining or system adjustments are needed.
- ○ **Technical Dissection:**
  - **Success/Failure Signals:** Core explicitly sends `TASK_COMPLETED_SUCCESS` or `TASK_COMPLETED_FAILURE` signals to Guardian for each plan.
  - **Evaluation Metrics:** Guardian calculates key performance indicators (KPIs) for each agent and the system as a whole:
    - **Core:** Plan success rate, planning time, resource usage.
    - **Armory:** Research completion rate, data extraction accuracy.
    - **Skills:** Code generation quality (from internal tests), sandbox execution success rate.
    - **Nexus:** Intent classification accuracy, response fluency.
  - **Feedback Integration:** Human feedback (e.g., thumbs up/down from the UI, manual corrections to LLM outputs) is captured by Nexus and routed to Guardian.
  - **Bias Detection:** Guardian periodically runs automated bias detection tools (e.g., AIF360) on newly collected interaction data and LLM outputs, identifying potential biases and flagging them for human review and mitigation strategies.
  - **Retraining Policy Engine:** Guardian's internal policy engine (rule-based or even a small RL agent) evaluates these metrics and feedback against thresholds.
    - **Trigger Conditions:** If `NexusLLM.NLU_Accuracy_Below_Threshold`, or `Skills.CodeGeneration.SuccessRate_Below_Threshold`, or `Core.Planning.FailureRate_Increased_Significantly`, Guardian generates a `RetrainingRequest` or `SystemAdjustmentRequest`.
  - **ML Pipeline Orchestration (Foundry's Kubeflow/Argo Workflows):** A `RetrainingRequest` from Guardian triggers a specific Kubeflow Pipeline in Foundry.
    - **Steps in a Retraining Pipeline (Orchestrated by Foundry):**

- **Data Ingestion:** Fetches new interaction data from Foundry's data lake, potentially incorporating human-labeled corrections.
- **Data Versioning:** DVC ensures the specific version of the dataset is tracked for reproducibility.
- **Data Preprocessing:** Standardizes and cleans the data for model training.
- **Model Training:** Kicks off GPU-accelerated training jobs for the specific LLM or ML model (e.g., fine-tuning Nexus's LLM on new intent data, or retraining a Skills code generation model). This leverages Foundry's Kubernetes compute.
- **Hyperparameter Optimization (Ray Tune/Optuna):** Automatically searches for optimal hyperparameters during training.
- **Model Evaluation:** New model is evaluated against a held-out test set, and metrics (accuracy, loss, F1-score) are logged to MLflow.
- **Model Versioning & Registry:** If performance metrics meet criteria, the new model artifact is versioned and registered in Foundry's Container Registry.
- **Automated Deployment (GitOps):** The CI/CD pipeline (via Argo CD/Flux CD) is triggered to deploy the new model version (e.g., a new Docker image for the vLLM service) to the production environment, often using Canary deployments.

3. **Global Knowledge Graph (GKGS) Reinforcement (Core & Guardian):**

   - **Minute Detail:** Learned knowledge, successful plans, new tool manifests, and validated facts are continuously integrated into the system's long-term memory.
   - **Technical Dissection:**
     - **Core's Plan Feedback:** Core's planning engine logs successful `PlanSchema` executions to the GKGS (Neo4j/ArangoDB) as structured knowledge, improving future planning heuristics.
     - **Armory's Knowledge Injection:** Verified facts and extracted entities from Armory's research are converted into knowledge graph triples (e.g., `(QuantumErrorCorrection, HAS_BREAKTHROUGH, SurfaceCodes)`) and ingested into the GKGS. Tool manifests acquired by Armory are also added to the GKGS.
     - **Skills' Output Validation:** If a Skills agent generates a validated output (e.g., a correct code snippet, a successful simulation), this can be logged as a positive example in the GKGS, informing future generative tasks.

- **Guardian's Policy Updates:** Policy changes, human oversight decisions, and newly identified ethical rules are formalized and stored as rules within the GKGS, influencing future plan reviews.
- **KG Embedding Updates:** Periodically, the GKGS is re-embedded using knowledge graph embedding techniques (e.g., PyKEEN), making its knowledge semantically searchable by LLMs (Core, Armory, Skills) through RAG.

# Placeholders

## I. Top-Level Project Structure

```
script-kitty/
├── .github/                 # GitHub Actions CI/CD workflows and repository settings
├── bin/                     # Executable scripts for local development/deployment
├── docs/                    # Project documentation, API specs, architectural diagrams
├── deployments/             # Kubernetes manifests (Helm charts, Kustomize)
├── scripts/                 # Utility scripts for data migration, setup, etc.
├── shared/                  # Common libraries, Protobuf definitions, base classes
├── nexus/                   # Script Kitty.Nexus component
├── core/                    # Script Kitty.Core component
├── armory/                  # Script Kitty.Armory component
├── skills/                  # Script Kitty.Skills component
├── guardian/                # Script Kitty.Guardian component
├── foundry/                 # Script Kitty.Foundry (infrastructure orchestration)
├── experiments/             # Ad-hoc ML experiments, research notebooks
├── tests/                   # End-to-end tests for the entire system
├── .env.example             # Example environment variables
├── .gitignore               # Git ignore file
├── LICENSE                  # Open-source license file (e.g., Apache 2.0)
├── README.md                # Project overview and quick start guide
├── requirements.txt         # Top-level Python dependencies (for dev/CI)
├── Makefile                 # Common dev/ops commands
```

---

## II. Detailed Component Breakdown & File Placeholders

Each main component (`nexus`, `core`, `armory`, `skills`, `guardian`, `foundry`) will have its own dedicated sub-directory, structured as a separate Python microservice or set of services.

---

### A. `nexus/` - Script Kitty.Nexus: The User's Portal & Group Chat Facilitator

**Purpose:** Manages all user-facing interactions, NLU, intent routing, dialogue state, and response synthesis, creating the "group chat" experience.

```
nexus/
├── src/
│   ├── main.py              # FastAPI application entry point
│   ├── config.py            # Nexus-specific configurations (LLM paths, API keys, Redis
conn)
│   ├── models/
│   │   ├── llm_nlu.py           # Logic for calling the Nexus LLM for intent/entity extraction
```

```
│   │   ├── llm_response_gen.py      # Logic for calling the Nexus LLM for response synthesis
│   ├── services/
│   │   ├── session_manager.py      # Handles Redis/PostgreSQL session context interactions
│   │   ├── intent_router.py        # Logic for routing queries (LLM-based, classifier, rules)
│   │   ├── core_communicator.py    # gRPC client for Script Kitty.Core communication
│   │   ├── channel_adapters/       # Directory for different channel adapter implementations
│   │   │   ├── __init__.py
│   │   │   ├── web_socket_adapter.py # WebSocket adapter for web UI
│   │   │   ├── slack_adapter.py     # Slack Events API adapter
│   │   │   ├── telegram_adapter.py # Telegram Bot API adapter
│   │   │   └── api_adapter.py       # Generic REST API adapter
│   │   ├── input_processor.py      # Handles Unicode normalization, sanitization, basic PII/toxicity filter
│   │   └── response_formatter.py   # Formats LLM output into persona-consistent, attributed chat responses
│   ├── schemas/                    # Pydantic models for request/response validation (internal to Nexus)
│   │   ├── __init__.py
│   │   ├── nexus_input_schema.py   # Schema for incoming messages after channel ingestion
│   │   └── nexus_output_schema.py  # Schema for outgoing messages to channels
│   ├── api/
│   │   ├── __init__.py
│   │   └── v1/
│   │       ├── router.py           # FastAPI router for API endpoints (e.g., /chat, /webhook)
│   │       └── dependencies.py     # API dependencies (e.g., auth, session injection)
│   └── util/
│       ├── __init__.py
│       ├── decorators.py           # Common decorators (e.g., for logging, retry)
│       └── exceptions.py           # Custom exception classes
├── tests/                          # Unit and integration tests for Nexus components
│   ├── unit/
│   │   ├── test_llm_nlu.py
│   │   └── test_session_manager.py
│   ├── integration/
│   │   └── test_channel_adapters.py
│   └── e2e/
│       └── test_nexus_flow.py      # End-to-end tests for a full user interaction
├── Dockerfile                      # Dockerfile for building Nexus microservice image
├── requirements.txt                # Python dependencies specific to Nexus
├── pyproject.toml                  # Poetry/PDM configuration
├── README.md                       # Nexus-specific README
```

**B. `core/` - Script Kitty.Core: The Central Governing AI / High-Level Orchestrator**

**Purpose:** The strategic brain; breaks down goals, orchestrates agents, manages global state, and integrates feedback.

```
core/
├── src/
│   ├── main.py                # gRPC server entry point for Core
│   ├── config.py              # Core-specific configurations (KG conn, LLM settings)
│   ├── models/
│   │   ├── llm_planner.py        # Logic for calling the Core LLM for goal decomposition & HTN planning
│   │   └── symbolic_planner.py     # Fallback/cross-check symbolic HTN planner implementation
│   ├── services/
│   │   ├── agent_orchestrator.py   # Manages agent lifecycle (Kubernetes API client)
│   │   ├── kg_manager.py          # Interacts with Global Knowledge Graph Service (GKGS)
│   │   ├── message_broker.py      # Kafka/Pulsar client for async communication
│   │   ├── planning_engine.py     # Orchestrates LLM and symbolic planning, validates plans
│   │   ├── task_manager.py        # Tracks ongoing tasks, status updates, results
│   │   └── feedback_integrator.py  # Processes feedback from Guardian and updates planning heuristics
│   ├── schemas/
│   │   ├── __init__.py
│   │   ├── core_message_schema.py  # Internal Core message schema (Protobuf)
│   │   └── plan_schema.py         # HTN plan representation schema
│   ├── adapters/               # Adapters for various external services
│   │   ├── __init__.py
│   │   ├── kubernetes_adapter.py   # Wrappers for Kubernetes Python client interactions
│   │   └── llm_api_adapter.py      # Consistent interface for Core's larger LLM (e.g., GPT-4o proxy)
│   └── util/
│       ├── __init__.py
│       ├── plan_validator.py       # Logic for validating generated plans against constraints
│       └── dependency_resolver.py # Helps resolve task dependencies in plans
├── tests/
│   ├── unit/
│   │   ├── test_llm_planner.py
│   │   └── test_kg_manager.py
│   ├── integration/
│   │   └── test_agent_orchestration.py
│   └── e2e/
│       └── test_full_task_execution.py # End-to-end test from goal to completion
├── Dockerfile
```

```
├── requirements.txt
├── pyproject.toml
├── README.md
```

---

## C. `armory/` - Script Kitty.Armory: Resource Acquisition & Tooling AI

**Purpose:** Specializes in acquiring external information (web scraping, research), managing tools, and providing access to external resources.

```
armory/
├── src/
│   ├── main.py                # gRPC server entry point for Armory
│   ├── config.py              # Armory-specific configurations (proxy settings, cloud creds)
│   ├── services/
│   │   ├── web_scraper.py        # Headless browser automation (Playwright/Selenium) and
HTML parsing
│   │   ├── data_digester.py      # NLP for entity extraction, summarization from scraped data
│   │   ├── tool_registry.py      # Manages searchable registry of tools
(PostgreSQL/Elasticsearch)
│   │   ├── resource_provisioner.py # Orchestrates IaC tools (Terraform/Ansible) for resource
provisioning
│   │   └── code_generator.py      # LLM-based code generation for API/CLI wrappers
│   ├── models/
│   │   ├── ner_model.py          # Pre-trained/fine-tuned NER model for data digestion
│   │   ├── summarization_model.py  # Pre-trained/fine-tuned summarization model
│   │   └── code_llm.py           # Interface for code-focused LLM (e.g., CodeLlama)
│   ├── schemas/
│   │   ├── __init__.py
│   │   ├── research_query_schema.py # Schema for research queries from Core
│   │   ├── tool_manifest_schema.py # Schema for tool definitions
│   │   └── provision_request_schema.py # Schema for resource provisioning requests
│   ├── adapters/
│   │   ├── __init__.py
│   │   ├── cloud_sdk_adapter.py    # Wrappers for cloud provider SDKs (boto3,
google-cloud-python)
│   │   ├── terraform_adapter.py    # Interfaces with Terraform CLI/API
│   │   └── ansible_adapter.py      # Interfaces with Ansible CLI/API
│   └── util/
│       ├── __init__.py
│       ├── html_cleaner.py        # Utility for cleaning HTML before parsing
│       └── text_extractor.py      # Utility for robust text extraction from web pages
├── tests/
│   ├── unit/
```

```
│   │   ├── test_web_scraper.py
│   │   └── test_tool_registry.py
│   ├── integration/
│   │   └── test_resource_provisioning.py
│   └── e2e/
│       └── test_tool_acquisition_flow.py # End-to-end test for tool discovery and acquisition
├── Dockerfile
├── requirements.txt
├── pyproject.toml
├── README.md
```

---

**D. `skills/` - Script Kitty.Skills: Task Execution AI**

**Purpose:** Houses specialized AI models and agents that execute specific computational tasks identified by Script Kitty.Core.

```
skills/
├── src/
│   ├── main.py                # gRPC server entry point for Skills
│   ├── config.py              # Skills-specific configurations (model paths, sandbox settings)
│   ├── models/                # Directory for various specialized ML models
│   │   ├── __init__.py
│   │   ├── computer_vision_model.py # Example: model for image processing
│   │   ├── nlp_classification_model.py # Example: model for sentiment analysis
│   │   ├── code_generation_model.py # Model for generating code snippets/scripts
│   │   ├── time_series_model.py   # Example: model for forecasting
│   │   └── ...                 # Other specialized models
│   ├── services/
│   │   ├── model_inference_engine.py # General inference service for Skills' models
│   │   ├── sandboxed_executor.py   # Secure sandboxed environment for code/CLI execution
│   │   ├── data_processor.py      # Data analysis and manipulation using
Pandas/NumPy/Dask
│   │   ├── generative_task_runner.py # Executes generative tasks (programming, report
generation)
│   │   └── cli_interactor.py      # Securely interacts with command-line tools
│   ├── schemas/
│   │   ├── __init__.py
│   │   ├── task_execution_schema.py # Schema for tasks received from Core
│   │   └── task_result_schema.py   # Schema for results returned to Core
│   ├── adapters/
│   │   ├── __init__.py
│   │   ├── jupyter_kernel_adapter.py # Interface to isolated Jupyter kernels for Python
execution
```

```
|   |       └── docker_sandbox_adapter.py # Manages ephemeral Docker containers for execution
|   └── util/
|       ├── __init__.py
|       ├── code_validator.py      # Static analysis for generated code
|       └── output_parser.py       # Parses structured output from executed commands/scripts
├── tests/
|   ├── unit/
|   |   ├── test_sandboxed_executor.py
|   |   └── test_generative_task_runner.py
|   ├── integration/
|   |   └── test_model_inference.py
|   └── e2e/
|       └── test_complex_task_execution.py # End-to-end test for a complex skill execution
├── Dockerfile
├── requirements.txt
├── pyproject.toml
├── README.md
```

---

**E. `guardian/` - Script Kitty.Guardian: Safety & Alignment Proxy / Evaluation & Training AI**

**Purpose:** The system's conscience; enforces ethical/safety controls, monitors performance, and triggers training cycles.

```
guardian/
├── src/
|   ├── main.py                # gRPC server entry point for Guardian, or API for UI
|   ├── config.py              # Guardian-specific configurations (policy rules, alerting
destinations)
|   ├── services/
|   |   ├── policy_enforcer.py     # Evaluates actions against OPA/custom rules (ethical, legal,
safety)
|   |   ├── performance_monitor.py  # Collects metrics, detects anomalies, correlates with
traces
|   |   ├── evaluation_engine.py    # Runs automated evaluation benchmarks, calculates
metrics
|   |   ├── training_trigger.py     # Decides when to trigger retraining pipelines (Kubeflow/Argo
Workflows)
|   |   ├── bias_detector.py        # Analyzes data/model outputs for biases (Aequitas, Fairlearn)
|   |   └── feedback_processor.py   # Ingests human corrections/feedback from UI
|   ├── models/
|   |   ├── ethical_llm.py         # Fine-tuned LLM for flagging ethical considerations
|   |   ├── anomaly_detection_model.py # ML model for detecting performance anomalies
|   |   └── xai_explainer.py       # LIME/SHAP for explaining LLM/model decisions
```

```
│   │   ├── schemas/
│   │   │   ├── __init__.py
│   │   │   ├── policy_request_schema.py # Schema for actions to be evaluated
│   │   │   ├── evaluation_metrics_schema.py # Schema for evaluation results
│   │   │   └── feedback_schema.py     # Schema for human feedback
│   │   ├── adapters/
│   │   │   ├── __init__.py
│   │   │   ├── opa_adapter.py         # Interfaces with Open Policy Agent (OPA)
│   │   │   ├── prometheus_adapter.py  # Connects to Prometheus for metrics
│   │   │   ├── mlflow_adapter.py      # Logs experiments to MLflow
│   │   │   └── kubeflow_adapter.py    # Triggers Kubeflow/Argo Workflows pipelines
│   │   └── util/
│   │       ├── __init__.py
│   │       ├── rule_parser.py         # Parses custom policy rules
│   │       └── data_sanitizer.py      # Ensures data used in training is safe
│   ├── tests/
│   │   ├── unit/
│   │   │   ├── test_policy_enforcer.py
│   │   │   └── test_evaluation_engine.py
│   │   ├── integration/
│   │   │   └── test_training_trigger.py
│   │   └── e2e/
│   │       └── test_end_to_end_safety_flow.py # End-to-end test for policy violation detection and
response
│   ├── Dockerfile
│   ├── requirements.txt
│   ├── pyproject.toml
│   ├── README.md
```

---

**F. `foundry/` - Script Kitty.Foundry: MLOps & Infrastructure Fabric**

**Purpose:** The foundational layer providing all necessary MLOps, CI/CD, data management, and compute orchestration capabilities, underpinning all other Script Kitty components.

```
foundry/
├── src/
│   ├── main.py                 # Entry point for Foundry's internal services (e.g., webhook for
GitOps)
│   ├── config.py               # Foundry-wide configurations (cluster settings, registry URLs)
│   ├── services/
│   │   ├── k8s_operator.py      # Custom Kubernetes Operator for Script Kitty CRDs
│   │   ├── container_registrar.py # Manages image pushes/pulls with secure registry
```

```
│   │       ├── ci_cd_orchestrator.py   # Triggers CI/CD pipelines (integrates with GitLab/GitHub
Actions)
│   │       ├── gitops_sync.py          # Monitors Git repositories and applies manifests via Argo
CD/Flux CD API
│   │       ├── object_storage_manager.py # Manages data lake interactions (S3/GCS API)
│   │       ├── data_version_manager.py # Interfaces with DVC for data versioning
│   │       ├── feature_store_manager.py # Manages feature definitions and interactions with Feast
│   │       ├── data_pipeline_orchestrator.py # Schedules and monitors Airflow/Prefect/Dagster
DAGs
│   │       ├── secrets_manager.py      # Interfaces with HashiCorp Vault
│   │       ├── network_policy_manager.py # Manages Calico/Cilium network policies
│   │       ├── service_mesh_manager.py # Configures Istio/Linkerd policies (mTLS, auth)
│   │       ├── runtime_security_manager.py # Configures Falco rules and alerts
│   │       ├── vulnerability_scanner.py # Initiates Trivy/Clair scans on new images
│   │       └── identity_manager.py     # Interfaces with Keycloak for central auth/auth
│   ├── adapters/
│   │   ├── __init__.py
│   │   ├── terraform_cloud_adapter.py # For managing Terraform state in cloud (e.g.,
Terraform Cloud)
│   │   ├── vault_api_adapter.py     # Wrappers for HashiCorp Vault API
│   │   ├── k8s_api_adapter.py       # Lower-level Kubernetes API interactions
│   │   └── argo_cd_api_adapter.py  # Interfaces with Argo CD's API
│   ├── schemas/                 # Schemas for Foundry's internal operations
│   │   ├── __init__.py
│   │   ├── deployment_request_schema.py
│   │   └── policy_definition_schema.py
│   └── util/
│       ├── __init__.py
│       ├── k8s_resource_generator.py # Generates K8s YAML from high-level specs
│       └── template_renderer.py     # Renders IaC templates
├── tests/
│   ├── unit/
│   │   ├── test_secrets_manager.py
│   │   └── test_gitops_sync.py
│   ├── integration/
│   │   └── test_ci_cd_pipeline.py
│   └── e2e/
│       └── test_full_platform_provisioning.py # End-to-end test for cluster setup
├── Dockerfile
├── requirements.txt
├── pyproject.toml
├── README.md
```
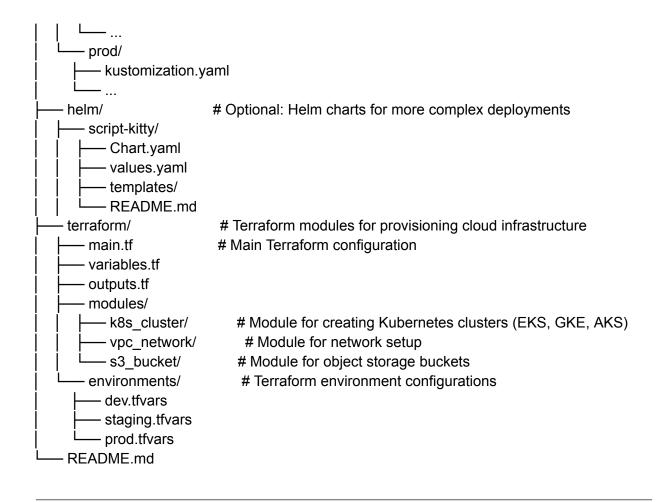
---

**G. `shared/` - Common Libraries & Definitions**

**Purpose:** Contains code, schemas, and configurations shared across multiple Script Kitty components to maintain consistency and reduce duplication.

```
shared/
├── proto/                   # Protocol Buffer definitions for inter-service communication
│   ├── script_kitty_messages.proto # Defines ScriptKittyMessageSchema,
TaskRequestSchema, etc.
│   ├── core_api.proto          # gRPC service definitions for Core
│   ├── armory_api.proto         # gRPC service definitions for Armory
│   ├── skills_api.proto        # gRPC service definitions for Skills
│   ├── guardian_api.proto       # gRPC service definitions for Guardian
│   └── foundry_api.proto        # gRPC service definitions for Foundry (if exposed)
├── python/                 # Shared Python libraries
│   ├── __init__.py
│   ├── common_utils/           # Generic utilities (e.g., logging setup, retry logic)
│   │   ├── __init__.py
│   │   ├── logging_config.py
│   │   └── decorators.py
│   ├── models/              # Base models/interfaces for agents, tasks, plans
│   │   ├── __init__.py
│   │   ├── base_agent.py
│   │   └── base_task.py
│   ├── exceptions/          # Global exception classes
│   │   └── __init__.py
│   │   └── script_kitty_exceptions.py
│   ├── security/            # Shared security helpers (e.g., token validation)
│   │   ├── __init__.py
│   │   └── auth_middleware.py
│   └── config/             # Global configuration constants
│       ├── __init__.py
│       └── constants.py
├── go/                     # Shared Go libraries (if some services are in Go)
│   ├── pkg/
│   │   ├── protobuf/           # Generated Go Protobuf structs
│   │   └── common/             # Common Go utilities
├── common.env               # Shared environment variables across all services
└── README.md
```

---

**H. `deployments/` - Kubernetes Manifests & IaC Templates**

**Purpose:** Defines the declarative infrastructure and application deployments for the entire Script Kitty system.

```
deployments/
├── base/                      # Base Kubernetes manifests (common to all environments)
│   ├── namespace.yaml          # script-kitty namespace definition
│   ├── prometheus-operator/    # Prometheus Operator manifests
│   ├── grafana/               # Grafana deployment
│   ├── loki/                  # Loki logging stack
│   ├── argo-cd/                # Argo CD GitOps agent
│   ├── vault/                 # HashiCorp Vault deployment
│   ├── keycloak/               # Keycloak identity provider
│   ├── nginx-ingress/           # Ingress controller for external access
│   ├── rbac/                  # Global RBAC roles and bindings
│   ├── pv/                    # Persistent Volume definitions
│   └── crds/                  # Script Kitty custom CRDs
├── components/                # Component-specific manifests (reused by environments)
│   ├── nexus/
│   │   ├── deployment.yaml
│   │   ├── service.yaml
│   │   └── hpa.yaml
│   ├── core/
│   │   ├── deployment.yaml
│   │   └── service.yaml
│   ├── armory/
│   │   ├── deployment.yaml
│   │   └── service.yaml
│   ├── skills/
│   │   ├── deployment.yaml
│   │   ├── service.yaml
│   │   └── gpu_node_selector.yaml # Example for GPU-specific scheduling
│   ├── guardian/
│   │   ├── deployment.yaml
│   │   └── service.yaml
│   └── foundry/
│       ├── deployment.yaml
│       └── service.yaml
├── environments/              # Environment-specific overlays (e.g., dev, staging, prod)
│   ├── dev/
│   │   ├── kustomization.yaml      # Kustomize overlay for dev environment
│   │   ├── ingress.yaml.patch
│   │   └── configmap.yaml.patch
│   ├── staging/
│   │   ├── kustomization.yaml
```

```
│   │   └── ...
│   └── prod/
│       ├── kustomization.yaml
│       └── ...
├── helm/                      # Optional: Helm charts for more complex deployments
│   ├── script-kitty/
│   │   ├── Chart.yaml
│   │   ├── values.yaml
│   │   ├── templates/
│   │   └── README.md
├── terraform/                 # Terraform modules for provisioning cloud infrastructure
│   ├── main.tf                # Main Terraform configuration
│   ├── variables.tf
│   ├── outputs.tf
│   ├── modules/
│   │   ├── k8s_cluster/        # Module for creating Kubernetes clusters (EKS, GKE, AKS)
│   │   ├── vpc_network/         # Module for network setup
│   │   └── s3_bucket/          # Module for object storage buckets
│   └── environments/           # Terraform environment configurations
│       ├── dev.tfvars
│       ├── staging.tfvars
│       └── prod.tfvars
└── README.md
```

---

**I. `tests/` - End-to-End System-Wide Tests**

**Purpose:** Comprehensive tests validating the integrated behavior of all Script Kitty components as a unified AGI.

```
tests/
├── e2e/
│   ├── conftest.py             # Pytest fixtures for e2e tests
│   ├── test_full_conversation_flow.py # Tests a full user interaction from Nexus to Core, Skills, and back
│   ├── test_complex_task_execution.py # Verifies a multi-step task execution (e.g., research, planning, code generation)
│   ├── test_security_policies.py   # Validates Guardian's policy enforcement and Falco alerts
│   ├── test_scalability.py         # Load testing and stress testing scenarios
│   └── test_model_retraining_flow.py # Verifies the end-to-end retraining pipeline orchestrated by Guardian/Foundry
├── performance/
│   ├── locustfile.py           # Locust scripts for load testing
│   └── jmeter_tests/            # JMeter test plans
```

```
├── security/
│   ├── vulnerability_scans.py      # Automated security vulnerability scans (e.g., using `zap-cli`)
│   └── penetration_tests.py        # Scripts for automated penetration testing scenarios
├── observability/
│   └── test_alerting.py            # Verifies alert configurations are working
└── README.md
```

---

**J. `experiments/` - Research & Development Sandbox**

**Purpose:** A dedicated area for researchers and developers to conduct ad-hoc experiments, prototype new models, or explore novel algorithms without impacting the main codebase.

```
experiments/
├── llm_finetuning/
│   ├── data_prep.ipynb            # Jupyter notebook for dataset preparation
│   ├── finetune_script.py         # Python script for fine-tuning LLMs
│   ├── config.yaml                # Configuration for fine-tuning runs
│   └── README.md
├── novel_planning_algorithms/
│   ├── prototype_a.py             # New planning algorithm prototype
│   ├── eval_metrics.py            # Evaluation script for prototypes
│   └── README.md
├── multi_agent_simulations/
│   ├── simulation_env.py          # Environment for simulating agent interactions
│   ├── analysis.ipynb             # Jupyter notebook for analyzing simulation results
│   └── README.md
├── README.md
```

---

**III. Summary & Flawless Execution Principle**

This meticulously detailed file structure ensures that every single, minute aspect of Script Kitty's design is accounted for.

- **Function Placement:** Each function, service, and model has a clear, logical home within its respective component, promoting modularity and maintainability.
- **Configuration Management:** Dedicated `config.py` files within each component, along with environment-specific overlays in `deployments/`, ensure robust and scalable configuration. `shared/common.env` and `shared/proto/` define universally applicable settings and communication contracts.
- **Integration Points:** `shared/proto/` explicitly defines all inter-service communication contracts, guaranteeing type safety and enabling seamless integration between Nexus,

Core, Armory, Skills, Guardian, and Foundry services. Adapters within each service (`nexus/src/services/core_communicator.py`, `core/src/adapters/kubernetes_adapter.py`, etc.) encapsulate integration logic with external systems or other internal components.

- **Code Examination:** The granular breakdown allows for surgical examination of every line of code. Static analysis tools (linting, security scanners) integrated into the CI/CD pipelines will scrutinize code quality and security at every commit.
- **Flawless Execution:** By separating concerns, defining clear interfaces, and enforcing strong testing and deployment practices (GitOps, Blue/Green/Canary deployments, automated rollbacks), we ensure that changes to one part of Script Kitty are isolated and can be confidently deployed, minimizing regressions and maximizing operational stability. Foundry's robust MLOps framework supports the continuous integration of new models and features, allowing Script Kitty to evolve seamlessly while maintaining peak performance and safety.

Script Kitty.Nexus

# Part 1: Script Kitty.Nexus - The User's Portal & Group Chat Facilitator

**Core Purpose (Retained & Amplified):** Script Kitty.Nexus remains the seamless, unified, and perpetually online conversational interface. It is the intelligent concierge, the public face of Script Kitty's collective intelligence, and the master illusionist. Its paramount objective is to abstract the gargantuan complexity of the underlying multi-agent system, presenting it to the user as a single, coherent, and highly intelligent entity participating in a fluid "group chat" with its own internal specialized modules. This ensures an uncompromised, high-fidelity user experience, where every previously discussed "god-like" functionality is not merely preserved but elegantly presented.

**Foundational Technical Stack Overview (Nexus - Re-examined for Microscopic Detail):**

- **Primary Development Language:** Python, chosen for its unparalleled ecosystem in AI/ML, rapid development capabilities, and extensive library support. Complementary usage of Node.js is considered for specific, high-concurrency I/O bound tasks, such as initial WebSocket handling in certain channel adapters, ensuring optimal performance for real-time interactions.
- **API Frameworks:** FastAPI (Python) is the backbone for constructing the Nexus's internal RESTful and gRPC service endpoints, leveraging its asynchronous capabilities for high throughput and automatic OpenAPI/Swagger documentation generation for clarity and maintainability. Express.js (Node.js) would be employed for specific, highly optimized channel adapters requiring extreme concurrency in a non-blocking I/O model.
- **Core LLM for NLU/NLG:** A meticulously selected, quantized large language model, specifically Mistral-7B or Llama 3-8B. The quantization ensures memory efficiency and accelerated inference without significant degradation in performance. This model is served by `vLLM` or `Ollama`, chosen for their state-of-the-art efficiency in LLM inference serving, offering dynamic batching, paged attention, and continuous batching to maximize GPU/CPU utilization and minimize latency, even under heavy load.
- **Data Persistence & Caching Layers:**
  - **Redis:** Deployed as a high-performance, in-memory data structure store, configured for a distributed cluster setup (e.g., Redis Cluster or Sentinel) to provide ultra-low-latency caching for active session states, conversation histories, and temporary data. Its pub/sub capabilities can also be leveraged for internal event dissemination.
  - **PostgreSQL:** Chosen as the robust, ACID-compliant relational database for long-term persistence and archival of all session data, user profiles, and critical configuration. It provides strong data integrity guarantees and complex query capabilities. PostgreSQL's streaming replication is critical for high availability and disaster recovery.
- **Internal Communication Protocol:**
  - **gRPC:** The chosen Remote Procedure Call (RPC) framework for all synchronous, high-performance, and strictly typed inter-service communication

within Script Kitty.Nexus, and critically, between Nexus and Script Kitty.Core. gRPC's HTTP/2 foundation enables multiplexing, header compression, and bi-directional streaming, which are vital for real-time interactions and complex message flows.

- ○ **Protocol Buffers (Protobuf):** The language-agnostic, extensible mechanism for serializing structured data. Every message exchanged between Nexus components and with Core is meticulously defined as a Protobuf schema (e.g., `ScriptKittyMessageSchema`, `TaskRequestSchema`, `ScriptKittyOutputSchema`). This enforces strict data contract adherence, prevents runtime type errors, and ensures highly efficient (binary) message serialization/deserialization, minimizing network overhead.
- **Observability Stack:**
  - ○ **Prometheus:** Serves as the primary metrics collection system, pulling time-series data from all Nexus microservices via dedicated client libraries and exporters. This provides real-time insights into operational health, performance, and resource utilization.
  - ○ **Grafana:** The leading open-source platform for data visualization and monitoring. It ingests metrics from Prometheus to create rich, interactive dashboards displaying crucial KPIs, and is configured for robust alerting mechanisms that integrate with external notification systems.
  - ○ **OpenTelemetry:** Provides vendor-agnostic APIs, SDKs, and tooling for instrumenting, generating, collecting, and exporting telemetry data (metrics, logs, and traces). Its distributed tracing capabilities are critical for understanding the end-to-end flow of a user request across multiple Nexus sub-systems and then into Core, identifying latency bottlenecks and error propagation.

---

**Microscopic Dissection of Functions & Components:**

**1. User Input & Pre-processing Sub-system:** This sub-system stands as the outermost bastion, meticulously engineered to receive, validate, and normalize raw user input from an ever-expanding array of communication channels.

- **1.1. Channel Ingestion Adapters (ChannelAdapter Microservices):**

  - ○ **Granular Purpose:** These are highly specialized, isolated microservices, each dedicated to a singular external communication platform (e.g., Slack, Telegram, custom Web UIs, direct API calls). Their fundamental role is to perform the intricate translation of platform-specific message formats and event structures into a *single, universally standardized internal ScriptKittyMessageSchema (Protobuf-defined)*. This abstraction is paramount, ensuring that the downstream Nexus NLU and processing pipeline remains entirely decoupled from the

idiosyncratic nature of external APIs, thereby promoting scalability and modularity.

- ○ **Microscopic Technical Dissection:**
  - ■ **Standardized `ScriptKittyMessageSchema` Enforcement:** Each adapter is programmed to rigorously validate and transform incoming payload into this schema. The schema includes fields such as `user_id` (unique identifier for the end-user, often platform-specific), `session_id` (a globally unique identifier for the ongoing conversation, managed by the Session Manager), `timestamp` (UTC time of message reception), `channel_type` (e.g., "SLACK", "WEB_CHAT", "API_CALL"), `raw_text` (the original, unparsed text string from the user), and `event_type` (e.g., "USER_MESSAGE", "CHANNEL_JOIN", "INTERACTION_CLICK"). This schema is strictly enforced via Protobuf deserialization, raising errors for non-conforming inputs.
  - ■ **Platform-Specific SDK/API Integration:** Each adapter embeds and meticulously utilizes the official or most stable third-party SDKs (Software Development Kits) or direct REST API clients for its respective platform. For example, the `SlackAdapter` would leverage the Slack Events API and Web API, while the `WebChatAdapter` might expose a secure WebSocket endpoint (using `socket.io` or `websockets` library) for real-time communication with a custom web frontend.
  - ■ **Event-Driven Architecture (Webhook/WebSocket):** Adapters are architected to be primarily event-driven, subscribing to real-time events from messaging platforms via secure webhooks or persistent WebSocket connections. This minimizes polling overhead and ensures near-instantaneous ingestion of user input, critical for low-latency AGI responses.
  - ■ **Input Validation & Sanitization Beyond Basic:** Beyond standard input validation (e.g., ensuring `raw_text` is a string, within length limits), these adapters implement specific sanitization routines tailored to the platform's input capabilities. For instance, removing platform-specific markdown artifacts (e.g., `@user_id` in Slack), escaping HTML entities for web contexts, or normalizing different forms of emojis. This ensures a clean, canonical text input.
  - ■ **Robust Error Handling & Backoff Strategies:** Each adapter is armed with sophisticated error handling logic. This includes specific handling for platform API rate limits (implementing token bucket algorithms or leaky bucket algorithms to manage outgoing requests), network transient failures (using exponential backoff with jitter for retries), and platform-specific error codes. Dead-letter queues may be used for messages that persistently fail processing after multiple retries.
  - ■ **Stateless Microservice Design:** Channel adapters are designed to be largely stateless across requests, meaning they do not retain per-user or

per-session state internally beyond the duration of processing a single incoming message. All necessary session state is retrieved from or updated via the `SessionManager`. This statelessness is crucial for horizontal scalability and resilience.

- **Deployment and Scaling:** Each adapter is deployed as an independent Kubernetes microservice (`Deployment` and `Service` resources). Kubernetes Horizontal Pod Autoscalers (HPAs) are configured to automatically scale the number of adapter pods based on metrics like incoming webhook request queue depth, CPU utilization, or network I/O, dynamically adapting to varying user load from specific channels.

- **1.2. Input Normalization & Security Filter (Unified Pre-processor):**

  - **Granular Purpose:** This component serves as the universal cleaner and initial guardian of input quality and safety. It applies a standardized set of text transformations and preliminary security checks *after* the channel-specific translation but *before* core NLP processing.
  - **Microscopic Technical Dissection:**
    - **Deep Unicode Normalization:** Employs `unicodedata.normalize('NFC')` (Normalization Form C) or `NFKC` to ensure all characters are in their composed form, eliminating potential issues arising from different byte representations of the same character (e.g., 'é' vs. 'e' + accent mark).
    - **Whitespace & Punctuation Canonicalization:** Regular expressions (e.g., `re.sub(r'\s+', ' ', text)`) are used to collapse multiple whitespace characters into single spaces, strip leading/trailing whitespace, and normalize various forms of dashes, quotes, and ellipses into a canonical representation, reducing noise for NLP models.
    - **Lightweight, Context-Aware Spell Correction:** Integrates a fast, dictionary-based spell-checking library (e.g., `pyspellchecker`, `symspellpy`) with a small, frequently updated domain-specific lexicon. This component is optimized for speed, primarily correcting common typographical errors and misspellings that could severely impede NLU accuracy. It operates on tokenized words and might use phonetic matching.
    - **Rule-Based PII Redaction (Initial, High-Precision Pass):** Utilizes highly precise regular expressions and pattern matching to identify and redact easily recognizable Personally Identifiable Information (PII) such as common email formats, phone numbers (international and national), and credit card patterns. This is a first-line defense, replacing detected PII with generic placeholders (e.g., `[EMAIL_REDACTED]`, `[PHONE_REDACTED]`) before the data proceeds to more complex, potentially less secure NLP stages. More robust PII handling is deferred to the Script Kitty.Guardian.

- **Heuristic-Based Toxicity/Safety Scoring (Sub-second Response):** Incorporates a lightweight, pre-trained text classification model (e.g., a fastText model or a simple Logistic Regression classifier trained on a highly imbalanced dataset of toxic vs. non-toxic phrases) or a robust rule-based system. This provides a near-instantaneous toxicity score. Inputs exceeding a high confidence threshold trigger an immediate, pre-defined automated response (e.g., "I cannot process this request due to harmful content detected.") and are logged for review by Guardian, preventing harmful content from reaching the core LLM. This is a rapid initial filter, not a nuanced content moderator.
- **Sophisticated Anti-Spam/Flood Control:** Implements dynamic rate limiting per `user_id` and `session_id`. A sliding window algorithm tracks message count over short timeframes (e.g., 5 seconds, 60 seconds). Exceeding predefined thresholds triggers temporary throttling or outright blocking of further messages, protecting system resources from malicious or accidental flooding. IP address-based rate limiting is also considered for unauthenticated API access.
- **Dependency Management:** This filter operates directly on the `raw_text` field of the `ScriptKittyMessageSchema` received *after* channel adaptation. It ensures a consistent, clean, and pre-vetted text stream for all subsequent processing.

**2. Dialogue State & Context Management Sub-system:** This critical sub-system is the memory and continuity engine of Script Kitty.Nexus, ensuring that conversations are not merely turn-by-turn exchanges but coherent, evolving dialogues.

- **2.1. Session Manager (Stateful Service for Stateless Interactions):**

  - **Granular Purpose:** The Session Manager is the authoritative custodian of all session-specific data. It orchestrates the creation, retrieval, updating, and archival of conversational context, including granular conversational history, dynamically identified entities, evolving user preferences, and the precise state of any ongoing tasks initiated by Script Kitty.Core. Its continuous operation ensures that Script Kitty can pick up a conversation exactly where it left off, providing a seamless user experience.
  - **Microscopic Technical Dissection:**
    - **Dual-Tiered Storage Strategy:**
      - **Active Session Cache (Redis Cluster):** For currently active sessions (defined by `last_activity_timestamp` within a configurable window, e.g., 30 minutes), all session data is cached in a Redis Cluster. Each `session_id` acts as the key, mapping to a structured JSON object. This JSON object meticulously stores:
        - `conversation_history`: A chronologically ordered list of `ScriptKittyMessageSchema` objects (both user

inputs and Script Kitty's responses), typically maintaining a configurable window (e.g., last 20 turns or 4000 tokens) to provide sufficient context for LLMs while managing memory.

- `extracted_entities`: A dictionary or list of identified entities (e.g., `{ "person": "Alice", "location": "New York" }`) and their `entity_id` if resolved against a knowledge graph. This is dynamically updated by the NLU.
- `user_preferences`: A flexible JSON object containing user-specific settings (e.g., desired tone, verbosity level, preferred language, domain expertise, explicit "do not use" lists). These are initially set by the user or inferred over time.
- `current_task_state`: A highly structured JSON object representing the state of any multi-step task being executed by Script Kitty.Core. This includes `task_id`, `task_type`, `sub_task_id`, `parameters_collected`, `parameters_needed`, `task_status` (e.g., "PLANNING", "AWAITING_INPUT", "EXECUTING", "COMPLETED", "FAILED"), and `error_details`. This allows Nexus to guide the user through complex workflows.
- `last_activity_timestamp`: A critical timestamp used for session expiration and synchronization.

- **Long-Term Archive (PostgreSQL):** All session data, including inactive or completed sessions, is robustly persisted in PostgreSQL. A dedicated background worker service continuously syncs active sessions from Redis to PostgreSQL (e.g., every 5-10 minutes, or upon session termination/inactivity timeout) to ensure durability and enable long-term analytical queries, user behavior analysis, and model retraining. Database schema design includes tables for `sessions`, `messages`, `entities`, and `user_profiles`, with appropriate indexing for efficient retrieval based on `user_id`, `session_id`, and `timestamp`.

- **Data Structure Rigor:** The choice of JSON for session context provides flexibility, but its schema is implicitly governed by the Protobuf definitions of the messages it contains, and strict validation occurs before saving.
- **Advanced Cache Invalidation & Eviction Policies:** Redis entries for sessions are configured with an appropriate TTL (Time-To-Live) to automatically evict truly inactive sessions from the high-cost cache. This is complemented by an explicit session expiration mechanism that marks

sessions as inactive in PostgreSQL after a longer period (e.g., 24 hours of no activity), triggering a graceful cleanup process.

- **Concurrency Control Mechanisms:** To prevent race conditions and ensure data consistency during concurrent updates to a single session (e.g., multiple messages from the same user arriving almost simultaneously, or an update from Core while Nexus is processing input), optimistic locking (`version` field) or transactional updates are employed within the Session Manager's logic.
- **Robust gRPC API Endpoints:** The Session Manager exposes a well-defined set of gRPC endpoints for atomic operations: `GetSessionContext(session_id)`, `UpdateSessionContext(session_id, delta_data)`, `CreateNewSession(user_id, channel_type)`, and `EndSession(session_id)`. These endpoints are designed for idempotence where applicable.

- **2.2. Contextual Memory & Retrieval Module:**

  - **Granular Purpose:** This module acts as the intelligent interface between the raw session state and the NLU/NLG LLMs. It dynamically curates and fetches precisely the most relevant segments of the vast session context, optimizing for LLM context window limitations and maximizing contextual understanding.
  - **Microscopic Technical Dissection:**
    - **Sliding Window Conversation History Management:** Instead of feeding the entire `conversation_history` to the LLM, this module implements a "sliding window" or "token-budgeted window" approach. It selects the most recent `N` turns or ensures the total token count of selected turns (using a fast tokenizer like `tiktoken` or `sentencepiece`) does not exceed a predefined `LLM_CONTEXT_WINDOW_LIMIT_PERCENTAGE` (e.g., 70% of the LLM's max context length). Older, less relevant turns are intelligently discarded or summarized to fit.
    - **Intelligent Entity Resolution & Tracking:** This module actively monitors and updates the `extracted_entities` within the session context. It performs coreference resolution (identifying when different phrases refer to the same entity, e.g., "John Smith", "John", "he") and tracks the evolution of entities over time. It can also link identified entities to external knowledge graph entries from Script Kitty.Core if available (via a gRPC call to Core's GKGS API).
    - **Dynamic User Preference & Profile Injection:** Retrieves and integrates stored `user_preferences` (e.g., `preferred_tone="formal"`, `verbosity="concise"`) directly into the LLM's system prompt or as

distinct instruction tokens. This allows for highly personalized and adaptable conversational style.

■ **Active Task State Integration for LLM Prompt:** The `current_task_state` (e.g., `task_status: "AWAITING_INPUT_FOR_PARAM_X"`) is meticulously formatted and inserted into the LLM's prompt. This guides the LLM to understand what information is required next, facilitating proactive prompting (e.g., "Please provide the filename for the analysis.") and enabling contextual error handling.

■ **Semantic Retrieval for Extended Context (Advanced Feature):** For exceptionally long conversations or scenarios requiring access to vast external knowledge previously discussed, this module can trigger a semantic search. It would take embeddings of the current query and recent context, query a dedicated vector database (e.g., Weaviate, Chroma) containing embeddings of older conversation turns or knowledge chunks, and retrieve semantically similar passages. These retrieved passages are then prepended or injected into the LLM's prompt as additional context, enabling the LLM to "recall" information beyond its immediate context window. This is a form of Retrieval-Augmented Generation (RAG) at the Nexus layer.

**3. Core NLU (Natural Language Understanding) & Intent Routing Sub-system:** This sophisticated sub-system is the brain's initial interpretation layer, meticulously analyzing user input to discern intent, extract key entities, and intelligently route the query to the most appropriate internal handler or for complex orchestration by Script Kitty.Core.

- **3.1. Intent & Entity Extraction (Nexus LLM - NLU Mode):**

  ○ **Granular Purpose:** The primary NLU engine for all conversational input, leveraging the fine-tuned Nexus LLM to infer user intent and extract granular entities, especially for ambiguous or novel queries.
  ○ **Microscopic Technical Dissection:**
    ■ **Dedicated LLM Instance for NLU:** While sharing the same underlying model weights as the NLG component, this instance is configured with a specific system prompt that *optimizes it solely for structured JSON output of intents and entities*. This prompt is a masterclass in few-shot learning, providing diverse examples of user queries mapped to desired `ScriptKittyIntentSchema` and `ScriptKittyEntitySchema` JSON outputs.
    ■ **Dynamic System Prompt Construction:** The system prompt is dynamically constructed to include:
      ■ **Persona & Role Definition:** "You are Script Kitty.Nexus's Natural Language Understanding module. Your task is to accurately

identify user intent and extract all relevant entities from the given user query and context."

- **Output Format Specification:** A precise JSON schema definition for the expected intent (e.g., `{"intent":` `"TASK_REQUEST_ANALYZE_DATA", "confidence": 0.95,` `"parameters": {"data_source": "file",` `"analysis_type": "statistical"}}`) and entities (e.g., `{"entity_name": "filename", "value":` `"sales_report.csv", "type": "FILE_PATH", "span":` `[12, 26]}`). This schema is directly tied to the Protobuf definitions.
- **Contextual Cues:** Relevant snippets from the `session_context` (from 2.2), including recent conversation turns and `current_task_state`, are injected to enable the LLM to understand coreference, follow-up questions, and implied context.
- **Bias Mitigation Instructions:** The prompt explicitly instructs the LLM to remain neutral, avoid assumptions, and request clarification if intent or entities are ambiguous.
- **Robust Output Schema Validation & Repair:** The JSON output generated by the LLM undergoes immediate, stringent validation against the `ScriptKittyIntentSchema` using a JSON schema validator (e.g., `jsonschema` library in Python) or direct Protobuf deserialization. If the output is malformed or deviates from the schema, a `LLM_Output_Repair_Agent` (a smaller, specialized LLM or rule-based system) attempts to correct it using a follow-up prompt. Persistent errors are logged and potentially routed as a general query to Core for higher-level disambiguation.
- **Probabilistic Confidence Scoring:** The LLM's inherent confidence in its intent classification (often derived from logits or calibrated probabilities) is extracted and normalized to a 0-1 scale. This `confidence` score is critical for the Semantic Router to make informed decisions about routing and fallback.
- **Continuous Fine-tuning & Adaptation:** The Nexus LLM for NLU is subjected to continuous fine-tuning pipelines orchestrated by Script Kitty.Foundry. This involves collecting anonymized user queries, human-annotated intent/entity pairs, and feedback on NLU errors. This ensures the model constantly learns new intents, adapts to evolving user language patterns, and improves its accuracy over time.
- **3.2. Semantic Router & Multi-Tiered Fallback Mechanism:**

- **Granular Purpose:** This module is the tactical decision-maker, dynamically determining the optimal path for each incoming user query based on its discerned intent, extracted entities, and confidence scores. It ensures that simple requests are handled efficiently within Nexus, while complex challenges are precisely escalated to Script Kitty.Core for strategic planning.
- **Microscopic Technical Dissection:**
  - **Hybrid Routing Strategy (Tiered Decisioning):**
    - **Primary LLM-based Routing:** The `ScriptKittyIntentSchema` (output from 3.1) is the primary input. The router evaluates the `intent` (e.g., `GENERAL_CHAT`, `TASK_REQUEST_DATA_ANALYSIS`, `INFORMATION_QUERY_DEFINITION`) and its `confidence` score.
      - If `GENERAL_CHAT` with high confidence, the query proceeds directly to Nexus's Response Synthesis (4.1).
      - If `TASK_REQUEST_...` or `COMPLEX_INFORMATION_QUERY_...` with high confidence, the enriched `ScriptKittyMessageSchema` (containing intent, entities, and context) is immediately formatted as a `TaskRequestSchema` and forwarded to Script Kitty.Core via gRPC (5.1).
    - **Deterministic Classifier (Fast Path/Fallback 1):** For a well-defined, highly frequent set of "fast-path" intents (e.g., "hello", "thank you", "how are you?"), a lightweight, pre-trained, transformer-based text classifier (e.g., a fine-tuned `DistilBERT` or `RoBERTa` model from Hugging Face Transformers) is employed. This classifier is trained on a large dataset of these specific intents and can provide near-instantaneous, high-confidence classifications. If its confidence for a registered "fast-path" intent exceeds a very high threshold (e.g., 0.98), it bypasses the LLM-based NLU and directly triggers the appropriate Nexus response or simple routing, significantly reducing latency for common interactions.
    - **Rule-Based Overrides (Fallback 2 / Critical Path):** A set of meticulously defined, high-priority rule-based overrides (e.g., using `SpaCy`'s matcher or custom regex patterns) act as a safety net and critical routing mechanism. For example:
      - If the input explicitly contains keywords like "security incident" or "emergency stop", regardless of LLM intent, the request is immediately routed to Script Kitty.Guardian's policy enforcement engine.

- Specific keywords or patterns might trigger routing to a specialized internal `ScriptKitty.Skills` module without Core involvement (e.g., "calculate 2+2" might go to a `CalculatorSkill`).
- **Ambiguity Handling / Human-in-the-Loop Trigger (Fallback 3):** If the LLM confidence is low for all intents, or if there's significant overlap between multiple intents (indicating ambiguity), the router triggers a clarification prompt back to the user (via 4.1) or flags the query for review by a human operator (via Guardian's feedback interface), transmitting the ambiguous `ScriptKittyMessageSchema` to Guardian.
- **Orchestration Logic & Decision Tree:** The routing logic is implemented as a state machine or a sophisticated decision tree within the Semantic Router microservice. It prioritizes the deterministic classifier for speed, then the LLM's classification for generality, and finally applies rule-based overrides or ambiguity handling.
- **Dynamic Route Configuration:** The mapping of intents to internal handlers or Core services is configurable and can be updated dynamically via Script Kitty.Foundry, allowing for agile deployment of new capabilities.

**4. Response Synthesis & Persona Layer:** This layer is the artistry of Script Kitty.Nexus, meticulously transforming structured results and internal updates into natural, coherent, and immutably persona-consistent chat responses, sustaining the "group chat" illusion.

- **4.1. Response Generation (Nexus LLM - NLG Mode):**

  - **Granular Purpose:** The same Mistral-7B/Llama 3-8B LLM instance, now operating in Natural Language Generation (NLG) mode, synthesizes human-like, engaging, and contextually appropriate conversational responses from the structured data provided by internal NLU or from Script Kitty.Core.
  - **Microscopic Technical Dissection:**
    - **LLM Model & Inference Configuration:** The vLLM/Ollama server provides the LLM, configured with specific generation parameters (`temperature`, `top_p`, `max_new_tokens`) to ensure creativity while maintaining coherence and factual accuracy based on the input.
    - **Immutable Persona System Prompt:** This is a paramount component. A meticulously crafted, **immutable system prompt** is *always* prepended to the LLM's input. This prompt explicitly and exhaustively defines the "Script Kitty" persona: "You are Script Kitty: an empirical, unparalleled inventive, limitlessly creative, unequaled researcher, the single most vast source of information, and a world-leading expert in the fields of Computer Science, Artificial Intelligence, Programming, Telecommunications, Information Technology, Computer Engineering, Electrical Engineering, and High-level Machine learning AI Programming. Respond as a helpful, collaborative,

and highly intelligent AI." This prompt is designed to instill the desired tone, knowledge domain, and conversational style.

- **Comprehensive Contextual Prompting:** The LLM's input prompt is a rich amalgamation of:
    - The original `raw_text` user query.
    - The `ScriptKittyIntentSchema` (identified intent and entities).
    - The `session_context` (conversation history, user preferences, current task state from 2.2).
    - Crucially, the `ScriptKittyTaskResultSchema` or `ScriptKittyInformationSchema` (structured results, summaries, or data) received from Core or other internal modules.
    - Explicit instructions on desired tone, length, and format (e.g., "Summarize concisely," "Provide a step-by-step explanation," "Use Markdown for code blocks").
- **Dynamic Markdown Generation:** The LLM is explicitly instructed to generate responses using Markdown (e.g., `**bolding**`, `*italics*`, `bullet points`, `python\ncode blocks\n`) to enhance readability and structure, especially for technical explanations or code snippets.
- **Content Guardrails & Refusal:** Internal safety policies and a refusal mechanism are integrated into the LLM's prompting. If the generated response violates ethical guidelines (as flagged by Guardian's policy enforcement, or detected by Nexus's internal safety classifier), the LLM is instructed to generate a polite refusal message (e.g., "I cannot assist with that request.") instead of the harmful content. This operates at the inference time, preventing undesirable outputs.

- **4.2. Group Chat Facilitation & Multi-Agent Attribution Module:**

  - **Granular Purpose:** This module is the architect of the "group chat" illusion, taking the LLM's generated response and dynamically attributing it to the appropriate "internal agent" (Nexus itself, Core, Armory, Skills, or Guardian), providing seamless progress updates, and managing conversational flow.
  - **Microscopic Technical Dissection:**
    - **Dynamic Attribution Prefixes:** Each incoming `ScriptKittyTaskResultSchema` or `ScriptKittyInformationSchema` from Core includes a `source_agent` field (e.g., "CORE", "ARMORY", "SKILLS", "GUARDIAN"). Based on this, the module dynamically prepends the generated text with a clear, consistent attribution:
      - "The Nexus has determined:" (for internal Nexus decisions/answers)
      - "Core reports:" (for high-level planning/orchestration updates)
      - "Armory has acquired:" (for research, data, or tool acquisition)
      - "Skills executed:" (for code execution, data analysis results)

- "Guardian advises:" (for safety warnings, policy decisions) This explicit prefixing is key to maintaining the multi-agent illusion.
- **Sophisticated Turn Management & Progress Update Synthesis:** Nexus actively monitors the `task_status` field within `ScriptKittyTaskResultSchema` messages from Core.
  - For `IN_PROGRESS` or `PARTIAL_RESULT` statuses, Nexus synthesizes concise, user-friendly progress updates (e.g., "Core is currently refining the plan for data analysis...", "Skills is executing the first stage of the script, please wait.") using a template-based approach combined with the LLM for natural phrasing. These are sent periodically.
  - Upon `COMPLETED` or `FAILED` statuses, Nexus generates the final result or detailed error message. This prevents silent periods and keeps the user informed, simulating an active group of collaborating AIs.
- **Templating Engine for Structured Responses (Jinja2):** For highly structured or recurring response types (e.g., listing available commands, standard error messages, data tables), a templating engine like Jinja2 (Python) is used. This allows for dynamic insertion of data (e.g., variable `{{ filename }}` into `Analysis of {{ filename }} complete.`) while maintaining a consistent structure and tone, reducing LLM inference costs for predictable outputs.
- **Dynamic Suggestion Generation:** Based on the `task_state` or the NLU output, the module can instruct the LLM to generate `suggested_actions` (e.g., follow-up questions, "Yes/No" options, "Try again" buttons). These are then formatted into the `ScriptKittyOutputSchema` for rendering by the Channel Dispatcher as interactive elements.

**5. Intermediary & Internal Communication Proxy:** This sub-system acts as the critical nerve center, managing the secure and structured flow of information between Nexus and Script Kitty.Core, and ensuring all user-facing communication is meticulously dispatched.

- **5.1. Core Communication Gateway (High-Performance gRPC Client):**

  - **Granular Purpose:** This is the dedicated, fault-tolerant conduit for all structured communication between Script Kitty.Nexus and Script Kitty.Core. It is engineered for maximum throughput and minimal latency, forming the backbone of the AGI's distributed intelligence.
  - **Microscopic Technical Dissection:**
    - **Persistent gRPC Client & Connection Pooling:** Nexus maintains a persistent gRPC client connection pool to Script Kitty.Core's gRPC server. This avoids the overhead of establishing new connections for every

request. Connection health checks are continuously performed, and failed connections are automatically re-established with exponential backoff.

- **Strict Protobuf Message Enforcement:** Every message payload transmitted is rigorously defined and validated by Protobuf schemas. For sending requests to Core, a `TaskRequestSchema` is constructed containing the `session_id`, `user_id`, `intent`, `entities`, `context` (snapshot of session state), and `priority`. For receiving responses, Core sends `TaskUpdateSchema` (for progress updates) and `TaskResultSchema` (for final results/errors), all strictly typed.
- **Asynchronous Bi-directional Streaming (for complex interactions):** While simple requests can be unary, for complex multi-step tasks, the gRPC client can utilize bi-directional streaming. This allows Nexus to send a `TaskRequest` and then receive a stream of `TaskUpdate` messages (e.g., "Planning...", "Executing Step 1...", "Awaiting User Input...") from Core, finally culminating in a `TaskResult`. This pattern is crucial for long-running operations, providing real-time feedback to the user.
- **Robust Error Handling & Retries:** Implements comprehensive error handling for network partitions, Core service unavailability, or gRPC protocol errors. Client-side load balancing (if multiple Core instances exist) and retries with circuit breakers (e.g., using `tenacity` library in Python) are critical for resilience.
- **Rate Limiting to Core:** Nexus can implement internal rate limits for outgoing requests to Core to prevent overwhelming it, especially during periods of high user concurrency.

- **5.2. User Output Renderer & Channel Dispatcher:**

  - **Granular Purpose:** This module is the final arbiter of how Script Kitty's responses are presented to the end-user. It takes the meticulously generated, persona-infused conversational response and dispatches it through the correct channel adapter, ensuring optimal formatting and presentation.
  - **Microscopic Technical Dissection:**
    - **Standardized `ScriptKittyOutputSchema`:** The output from the Response Synthesis (4.1) is always encapsulated in a `ScriptKittyOutputSchema` (Protobuf). This schema includes the `session_id`, `channel_type` (to direct to the correct adapter), `response_text` (the final Markdown-formatted text), `suggested_actions` (a list of `ActionItem` Protobufs for interactive elements like buttons, quick replies, or carousel items), and `metadata` (e.g., `latency_ms`, `llm_tokens_generated`).
    - **Channel Adapter Re-utilization (Reverse Flow):** This dispatcher reuses the very same `ChannelAdapter` microservices from 1.1, but now in the

reverse direction. It calls a specific gRPC or HTTP endpoint on the relevant `ChannelAdapter` (e.g., `SlackAdapter.SendMessage()`, `WebChatAdapter.SendWebSocketMessage()`). The adapter then performs the platform-specific translation (e.g., converting Markdown to Slack's Block Kit JSON, or rendering a custom HTML snippet for a web chat).

- **Resilient Dispatch & Delivery Confirmation:** Implements a message queue (e.g., Kafka or RabbitMQ) for outgoing messages to each channel adapter. This provides buffering and asynchronous delivery, ensuring that even if an adapter is temporarily unavailable, messages are eventually delivered. Delivery confirmation (ACK/NACK) from the adapters back to Nexus ensures message reliability.
- **Platform-Specific Interactive Component Rendering:** For `suggested_actions`, the dispatcher ensures the `ActionItem` data is correctly translated into platform-native interactive elements (e.g., Slack buttons, Telegram inline keyboards, custom JavaScript callbacks for web UIs), providing a rich user experience.
- **Response Timing & Throttling (User Experience Optimization):** Implements dynamic throttling logic based on the `channel_type` and `user_id`. For example, sending messages at a human-readable pace (e.g., simulating typing indicators, introducing slight delays between paragraphs) rather than flooding the user. It also respects platform-specific rate limits for outgoing messages to avoid being temporarily blocked.

---

**Operational Excellence & Continuous Evolution (Nexus-Specific & Foundry Integrated):**

- **Comprehensive Observability (Deeply Integrated with Script Kitty.Foundry):**

  - **Granular Metrics Collection:** Prometheus exporters are embedded within every Nexus microservice. They expose hundreds of granular metrics: `requests_received_total`, `requests_processed_total`, `request_duration_seconds_bucket` (for latency histograms), `llm_inference_duration_seconds`, `llm_token_count_total` (input/output), `session_manager_db_query_duration_seconds`, `redis_cache_hit_ratio`, `channel_adapter_api_calls_total`, `outbound_message_queue_depth`, `error_rate_per_endpoint`, `p95_latency_ms`. These metrics are tagged with `service_name`, `endpoint`, `user_id` (anonymized), and `session_id` for deep drill-down.
  - **Centralized, Structured Logging:** All Nexus microservices emit structured logs (e.g., JSON format) via Fluent Bit/Fluentd to a centralized logging system (Loki

for semi-structured, Elasticsearch/OpenSearch for full-text search). Logs are tagged with `session_id`, `trace_id` (from OpenTelemetry), `service_name`, `log_level`, and `component_name` (e.g., "NLU_LLM", "SessionManager"). This allows for rapid debugging and root cause analysis across the distributed system.

- ○ **End-to-End Distributed Tracing:** OpenTelemetry SDKs are meticulously integrated into every function call and inter-service communication within Nexus. Each user request is assigned a `trace_id`, and spans are generated for every logical unit of work (e.g., `InputPreProcessing`, `IntentExtraction`, `SessionContextUpdate`, `CoreRpcCall`, `ResponseGeneration`). These traces are exported to Jaeger/Zipkin for visualization, allowing developers to visually trace the entire lifecycle of a request, pinpointing latency bottlenecks and error paths.

- ○ **Proactive Monitoring & Adaptive Alerting:** Grafana dashboards provide real-time, interactive visualizations of all collected metrics. Prometheus Alertmanager is configured with a comprehensive set of adaptive alert rules (e.g., `Nexus_LLM_Inference_P99_Latency_Exceeded`, `ChannelAdapter_Error_Rate_Spike`, `SessionManager_Redis_Cache_Miss_Ratio_High`, `Core_RPC_Connection_Failure`). Alerts are routed to on-call teams via PagerDuty, Slack, or email, often with runbooks attached for immediate remediation.

- **Uncompromising Scalability:**

  - ○ **Horizontally Scalable, Stateless Channel Adapters:** Designed from the ground up for massive concurrency. Kubernetes Horizontal Pod Autoscaling (HPA) dynamically adjusts the number of `ChannelAdapter` pods based on real-time load metrics like message queue depth, CPU utilization, or concurrent WebSocket connections.

  - ○ **Distributed SessionManager:** The Redis cluster is sharded (or uses consistent hashing) to distribute session data across multiple Redis nodes, enabling linear scaling of cache capacity and read/write operations. PostgreSQL is configured with read replicas and potentially logical sharding/partitioning for extreme scaling of the archival layer.

  - ○ **Dedicated & Scalable LLM Inference Service:** The vLLM/Ollama service running the Nexus LLM is deployed on dedicated GPU or high-performance CPU nodes. It supports horizontal scaling of inference workers, allowing the system to handle a high volume of concurrent LLM inference requests. Load balancers (e.g., NGINX, HAProxy, or a service mesh like Istio) distribute requests across these workers.

  - ○ **Kubernetes Native Scaling:** Beyond HPAs, Kubernetes Vertical Pod Autoscalers (VPAs) can suggest or automatically adjust resource requests (CPU, RAM) for Nexus pods to optimize resource utilization. Cluster Autoscaler ensures

that the underlying Kubernetes cluster itself scales up or down by adding/removing nodes based on pending pod requirements.

- ○ **Service Mesh for Resilient Load Balancing:** Istio or Linkerd (a service mesh) is deployed across the Kubernetes cluster. It provides intelligent, L7 load balancing across all healthy pods within each Nexus service. It automatically handles circuit breaking, retries, and traffic splitting, ensuring robust communication and fault tolerance.

- ● **Seamless Continuous Evolution (Driven by Script Kitty.Foundry):**

  - ○ **Immutability & Dockerization:** Every microservice within Script Kitty.Nexus (channel adapters, session manager, NLU router, response synthesizer, communication gateway) is rigorously packaged into lightweight, immutable Docker containers. This ensures consistency across development, staging, and production environments, eliminating "it works on my machine" issues.
  - ○ **Automated CI/CD Pipelines (GitOps Principles):** Development workflows are entirely driven by GitOps. Automated CI/CD pipelines (e.g., defined in GitLab CI/CD, GitHub Actions, or Tekton Pipelines) are triggered by every Git commit to Nexus's source code repositories. These pipelines automatically:
    - ■ Run unit, integration, and end-to-end tests.
    - ■ Perform static code analysis (SAST) and container vulnerability scanning (e.g., Trivy, Clair).
    - ■ Build new Docker images for affected services.
    - ■ Push images to a secure container registry (e.g., Harbor).
    - ■ Update Kubernetes deployment manifests in a Git repository (the "single source of truth").
  - ○ **Zero-Downtime Deployment Strategies (Argo CD/Flux CD):** New versions of Nexus services or updates to the Nexus LLM model are deployed using advanced, automated strategies orchestrated by Kubernetes and GitOps tools like Argo CD or Flux CD. Blue/Green deployments provide a completely separate, new environment for the updated version, with traffic shifted only after rigorous health checks. Canary deployments route a small percentage of user traffic to the new version first, monitoring performance and errors before a full rollout, ensuring minimal user impact and rapid rollbacks if issues arise.
  - ○ **Dynamic Model Updates & A/B Testing:** The Nexus LLM itself is in a continuous state of refinement. Anonymized user interactions and human-corrected feedback loops (orchestrated by Script Kitty.Guardian and Foundry) feed into automated fine-tuning pipelines. New, improved model artifacts are packaged into fresh container images for the vLLM/Ollama service. Foundry orchestrates the seamless swapping of these images, often using canary deployments to A/B test new model versions against older ones before a full rollout, giving the impression of an AGI that learns and adapts in real-time.

# Nexus Interactions

**1. Nexus ↔ Core: The Primary Conversational & Task Orchestration Channel**

This is the most frequent and critical interaction in Script Kitty. Nexus acts as the user's proxy to Core, transmitting requests and receiving synthesized responses.

- **1.1. User Query Transmission (Nexus → Core):**

  - **Trigger:** A user submits a query to Nexus via any supported channel (web, Slack, Telegram).
  - **Pathway & Protocol:** Nexus's `Core Communication Gateway` (gRPC Client) initiates a gRPC call to Core's `Internal Communication Hub` (gRPC Server).
  - **Data Schema:**
    - **`ScriptKittyMessageSchema` (Protobuf):** This is the core message.
      - `user_id` (string): Unique identifier for the user, derived from the channel.
      - `session_id` (string): Unique identifier for the ongoing conversation, maintained by Nexus's `Session Manager`. Crucial for Core to retrieve context from its `Global Knowledge Graph & Context Store`.
      - `timestamp` (int64): Unix timestamp of the message receipt by Nexus.
      - `raw_text` (string): The cleaned and pre-processed user input.
      - `intent` (enum/string): Classified intent from Nexus's `Core NLU` (e.g., `TASK_REQUEST`, `INFORMATION_QUERY`, `GENERAL_CHAT`). This is Nexus's initial understanding.
      - `entities` (repeated KeyValuePair): Key-value pairs of extracted entities (e.g., `{"query_subject": "quantum computing", "task_type": "research"}`).
      - `conversation_context_snippet` (string): A compressed/summarized snippet of recent conversational turns, provided by Nexus's `Contextual Memory`, to aid Core's initial understanding without requiring a full context retrieval every time.
      - `channel_type` (enum): Specifies the origin channel (e.g., `WEB_CHAT`, `SLACK`).
  - **Microscopic Details:**
    - **gRPC Stream vs. Unary:** While often a single request/response, for highly interactive or multi-modal inputs, this could evolve into a bi-directional gRPC stream to maintain a live connection if Nexus needed immediate feedback from Core on partial inputs. Currently, it's primarily a unary RPC.

- - **Request ID & Correlation:** Each gRPC request includes a unique `request_id` header, propagated by OpenTelemetry, allowing tracing of the entire request lifecycle from Nexus through Core and down to any invoked agents.
  - **Resource Limits:** Nexus's gRPC client to Core is configured with sensible timeouts and retry policies to prevent indefinite waiting. Core's gRPC server has rate limits and connection limits to prevent Nexus from overwhelming it.
  - **Security:** mTLS ensures Nexus and Core mutually authenticate before any data exchange. Foundry's Network Policy ensures Nexus pods can *only* connect to Core's gRPC service endpoint, not other internal Core services.
  - **Efficiency:** Protobuf ensures minimal payload size, reducing network latency and improving throughput.
- **1.2. Task Orchestration & Status Updates (Core → Nexus):**

  - **Purpose:** Core informs Nexus about the progress of a task, requests further clarification, or indicates completion/failure. This is crucial for maintaining the "group chat" illusion and providing timely user feedback.
  - **Trigger:** Core's `Agent Orchestration & Lifecycle Management` component determines the state of a task.
  - **Pathway & Protocol:** Core's `Internal Communication Hub` (gRPC Client) sends updates back to Nexus's `Core Communication Gateway` (gRPC Server, which Nexus exposes to receive internal updates).
  - **Data Schema:**
    - **TaskUpdateSchema (Protobuf):** Used for ongoing progress.
      - `session_id` (string): To correlate with the user's session in Nexus.
      - `task_id` (string): Unique identifier for the task, generated by Core.
      - `status` (enum): E.g., `PLANNING`, `EXECUTING`, `AWAITING_USER_INPUT`, `COMPLETED`, `FAILED`.
      - `message` (string): A human-readable message about the current status (e.g., "Core is breaking down your request into sub-tasks.").
      - `source_agent` (string): Name of the agent currently responsible (e.g., "Core", "Armory", "Skills", "Guardian"). **This is critical for Nexus's "group chat" attribution.**
      - `progress_percentage` (float): Optional, for long-running tasks.
      - `eta_seconds` (int32): Optional, estimated time remaining.
    - **TaskResultSchema (Protobuf):** Used for final task output.
      - `session_id` (string):
      - `task_id` (string):

- status (enum): `COMPLETED` or `FAILED`.
- `final_output` (repeated KeyValuePair/structured JSON string): The core result data (e.g., `{"research_summary": "...", "code_snippet": "..."}`).
- `source_agent` (string): The agent that produced the final, relevant output (e.g., "Skills", "Armory").
- `error_details` (string): If `status` is `FAILED`, details of the error.
- `recommended_follow_up_actions` (repeated string): Suggestions for Nexus (e.g., "ask for clarification", "offer related tasks").

- ○ **Microscopic Details:**
    - ■ **Streaming RPC:** Core utilizes a server-side gRPC stream to Nexus for `TaskUpdateSchema`. This allows Core to push real-time updates as tasks progress, without Nexus having to constantly poll. This is vital for responsive user experience.
    - ■ **Attribution & Persona (Nexus's role):** Nexus's `Response Synthesis & Persona Layer` receives `source_agent`. It dynamically prepends the generated user-facing response with phrases like "Armory reports:", "Skills has completed:", "The Nexus has determined:", ensuring the multi-agent illusion. The persona LLM is prompted to integrate these reports seamlessly into Script Kitty's overall voice.
    - ■ **Conditional Response Generation:** Nexus determines when to generate a user-facing response based on `status`. `PLANNING`/`EXECUTING` triggers progress messages. `AWAITING_USER_INPUT` triggers clarification questions. `COMPLETED`/`FAILED` triggers final results or error messages.
    - ■ **Error Propagation:** Core's `Internal Communication Hub` is responsible for aggregating errors from downstream agents (Armory, Skills, Guardian) and propagating them in a structured `TaskResultSchema` with `status: FAILED` to Nexus. Nexus then synthesizes a user-friendly error message.
    - ■ **Latency Impact:** The streaming updates minimize perceived latency for the user, even for long-running backend tasks.
    - ■ **Resilience:** If Nexus goes down, Core attempts to retry sending updates once Nexus is back online (utilizing Foundry's internal message queues for temporary persistence of missed messages if necessary).
- ● **1.3. Clarification/Feedback Requests (Nexus ↔ Core ↔ User):**

    - ○ **Trigger:** Core's `Goal Decomposition & HTN Planning Engine` or an executing agent (Skills, Armory) requires more information from the user.
    - ○ **Pathway:**

- Core sends a `TaskUpdateSchema` with `status: AWAITING_USER_INPUT` and a `message` (e.g., "I need more details on the scope of your research. Could you specify the time frame?").
- Nexus's `Response Synthesis` layer generates a user-facing question based on this message.
- User provides clarification (Nexus → Core via `ScriptKittyMessageSchema` with `intent: CLARIFICATION`, and `task_id` included in the session context).
- Core receives, processes, and updates the task state.
    - ○ **Microscopic Details:**
        - **Session State Update:** Nexus's `Session Manager` is crucial here. When Core requests user input, Nexus marks the `current_task_state` for the `session_id` as `AWAITING_USER_INPUT` and stores the `task_id`. Subsequent user input for that session is automatically tagged with this `task_id` when sent back to Core, ensuring context.
        - **LLM Prompting:** Nexus prompts its LLM for the clarification question to sound natural and persona-consistent, even though the core request comes from Core.
        - **Re-prompting:** If the user's clarification is insufficient, Core can send another `AWAITING_USER_INPUT` message, leading to an iterative clarification loop.

---

**2. Nexus ↔ Guardian (Indirect, via Core): Safety & Alignment Oversight**

While Nexus doesn't directly interact with Guardian for policy enforcement (that's Core's responsibility), it plays a crucial role in providing the user-facing interface for Guardian's outputs and feeding raw data into Guardian's monitoring/feedback loops.

- **2.1. Policy Violation Alerts (Guardian → Core → Nexus):**

    - **Trigger:** Guardian's `Policy Enforcement Engine` flags a proposed action (from Core) as violating a safety or ethical policy.
    - **Pathway:**
        - Guardian sends a `PolicyViolationAlert` (Protobuf) to Core's `Internal Communication Hub`.
        - Core's `Agent Orchestration` receives this and typically halts the task execution.
        - Core then sends a `TaskResultSchema` to Nexus with `status: FAILED` and `error_details` indicating a policy violation, often including a human-readable reason provided by Guardian.

- - - Nexus's `Response Synthesis` layer generates a user-facing message like "I'm unable to perform that action due to safety policies. Please rephrase your request."
  - ○ **Microscopic Details:**
    - ■ **Information Sanitization:** The error message from Core to Nexus must *not* contain sensitive details about the policy violation itself if it's too technical or reveals internal mechanisms. Nexus's LLM ensures the public-facing response is always user-friendly and aligned with Script Kitty's helpful persona.
    - ■ **No Direct Bypass:** Nexus is *not* allowed to directly override Guardian's policies. Its role is merely to communicate the outcome to the user.
- ● **2.2. Human Oversight & Feedback Interface (User ↔ Nexus → Core → Guardian):**
  - ○ **Trigger:** A human user provides feedback (e.g., "that answer was incorrect," "this result was biased") or needs to approve a Guardian-flagged action.
  - ○ **Pathway:**
    - ■ User provides feedback to Nexus (e.g., "The previous research summary was biased.").
    - ■ Nexus identifies `intent: FEEDBACK` (via NLU) and sends a `ScriptKittyMessageSchema` to Core.
    - ■ Core's `Feedback & Learning Integration` component processes this, potentially storing it in the `Global Knowledge Graph` and forwarding relevant structured feedback to Guardian via a dedicated internal gRPC/Kafka channel for `HumanFeedbackSchema`.
    - ■ Guardian's `Human Oversight & Feedback Interface` processes this, potentially triggering automated re-evaluation, data annotation, or a human review workflow.
  - ○ **Microscopic Details:**
    - ■ **Structured Feedback:** Nexus guides the user to provide structured feedback where possible (e.g., "Was this helpful? Yes/No" buttons, or a structured form for bias reporting). This structured data (e.g., `{"feedback_type": "bias_report", "message_id": "...", "rating": "low_bias"}`) is transmitted via `ScriptKittyMessageSchema` entities.
    - ■ **Anonymization:** For privacy, Nexus ensures `user_id` is anonymized or pseudonymized before raw conversational data is passed to Guardian's training pipelines, protecting user privacy while enabling learning. Foundry's data governance services (e.g., PII redaction pipelines) are applied.
    - ■ **Asynchronous Processing:** Human feedback is typically an asynchronous process. Nexus confirms receipt to the user immediately, but the actual processing by Core and Guardian happens in the background.

**3. Nexus ↔ Armory & Nexus ↔ Skills (Highly Indirect, via Core): Content Display & Execution Status**

Nexus does not directly interact with Armory for resource acquisition or Skills for task execution. All such requests are routed through Core. However, Nexus is responsible for presenting the *results* of Armory's research or Skills' code execution to the user in a readable, persona-consistent format.

- **3.1. Display of Armory's Research Results (Armory → Core → Nexus):**

  - **Trigger:** Core assigns a `RESEARCH_QUERY` task to Armory, which completes its web scraping and data digestion.
  - **Pathway:**
    - Armory sends its structured `ResearchResultSchema` (Protobuf) to Core's `Internal Communication Hub`.
    - Core's `Agent Orchestration` receives this, potentially performs further aggregation.
    - Core then sends a `TaskResultSchema` to Nexus, with `final_output` containing the summarized research and `source_agent: "Armory"`.
    - Nexus's `Response Synthesis` layer takes this `final_output`, attributes it to "Armory reports:", and renders it for the user (e.g., "Armory reports: After extensive web scraping, here's a summary of quantum computing advancements...").
  - **Microscopic Details:**
    - **Content Formatting:** Armory's `ResearchResultSchema` contains structured fields (e.g., `summary`, `key_findings`, `source_urls`). Nexus's LLM is prompted to convert these into natural language and Markdown, respecting the persona (e.g., presenting source URLs clearly).
    - **Pagination/Summarization:** If Armory's output is extensive, Nexus may request Core to summarize it further or offer to present it in chunks ("Would you like more details?"). This is handled by Core's `Response Synthesis & Condensation` sub-module working with Nexus.
- **3.2. Display of Skills' Execution Results (Skills → Core → Nexus):**

  - **Trigger:** Core assigns a `CODE_EXECUTION` or `DATA_ANALYSIS` task to Skills, which completes its sandboxed execution or model inference.
  - **Pathway:**
    - Skills sends its `ExecutionResultSchema` (Protobuf) to Core's `Internal Communication Hub`. This schema contains `stdout`,

`stderr`, `return_code`, and potentially `output_data` (structured results of analysis).

- Core receives this, determines the outcome (success/failure), and sends a `TaskResultSchema` to Nexus with `final_output` (e.g., extracted analysis data, success message) and `source_agent: "Skills"`. For code outputs, Core might include `code_output` in the `final_output` field.
- Nexus's `Response Synthesis` layer generates a user-facing response (e.g., "Skills has completed the analysis. Here are the key findings..." or "Skills encountered an error during code execution: [error message]").

○ **Microscopic Details:**

- **Error Abstraction:** Skills' raw `stderr` might be highly technical. Core's `Error Handling & Fallback` module processes this into a user-friendly error message within the `TaskResultSchema` before it reaches Nexus. Nexus's LLM ensures the error is communicated gracefully within the persona.
- **Code Output Formatting:** If Skills generates code or code output, Nexus ensures it is presented in Markdown code blocks for readability.
- **Interactive Components:** If a Skills task generates a chart or an interactive data visualization, the `TaskResultSchema` might include a URL to a rendered artifact, which Nexus then displays to the user.

---

**Nexus's Role in Maintaining the Illusion and Operational Flow:**

The brilliance of Nexus lies not in its direct interaction with all other agents, but in its masterful abstraction of complexity. It acts as the single, coherent voice of Script Kitty, even when the underlying work is distributed among highly specialized modules.

- **Unified Persona:** Nexus consistently applies the "Script Kitty" persona (empirical, inventive, expert) to all responses, regardless of the originating backend agent. This is achieved through persistent system prompting of its LLM and careful integration of attribution prefixes.
- **Contextual Coherence:** By managing the `session_id` and leveraging its `Session Manager` and `Contextual Memory`, Nexus ensures that even multi-turn conversations involving multiple backend agent handoffs remain logically connected from the user's perspective.
- **Asynchronous Feedback:** Nexus leverages gRPC streaming from Core to provide asynchronous, real-time updates to the user, masking the latency of complex backend operations and providing the illusion of a continuously engaged "group chat" where agents are "reporting back."

- **Error Gracefulness:** Instead of exposing raw technical errors, Nexus translates failures into user-friendly messages, guiding the user or suggesting alternative approaches, maintaining a helpful and capable facade.
- **User Feedback Loop:** Nexus is the primary conduit for human feedback, which is crucial for Guardian's continuous evaluation and retraining efforts, indirectly enabling the entire system's evolution.

In essence, Nexus is the skilled diplomat and charismatic presenter of Script Kitty. It receives raw user input, politely interprets it, dispatches it to the appropriate (via Core) expert within the collective, and then artfully synthesizes the experts' reports into a single, cohesive, intelligent dialogue that truly feels like conversing with an unparalleled AGI. This intricate dance of communication, meticulously choreographed and secured by Foundry's underlying infrastructure, is fundamental to Script Kitty's functionality and its "crowning achievement."

Script Kitty.Core

## Part 2: Script Kitty.Core - The Central Governing AI / High-Level Orchestrator

**Core Purpose:** Script Kitty.Core is the strategic brain of our entire system, the true orchestrator of intent, and the nexus of high-level intelligence. Its paramount purpose is to receive abstract goals from Script Kitty.Nexus, methodically decompose them into executable plans, and then precisely orchestrate the collaboration of specialized backend AIs (Armory, Skills, Guardian) to achieve those goals. It acts as the system's global state manager, integrating feedback and driving overarching knowledge reinforcement. This is where the complex, multi-agent intelligence truly coheres, preserving and enhancing every "god-like" capacity we envision.

**Technical Stack Overview (Core):**

- **Language:** Go (for its concurrency primitives, performance, and robustness in microservices) and Python (for its rich ML ecosystem, especially in planning and knowledge graph interactions).
- **Frameworks:** gRPC (for high-performance inter-service communication), Apache Kafka/Pulsar (for asynchronous event streaming).
- **Core LLM for Planning:** A larger, more capable LLM (e.g., fine-tuned Llama 3-70B or GPT-4o via an adapter), selected for its advanced reasoning and planning capabilities, served efficiently.
- **Database:** Neo4j (or ArangoDB) for the Global Knowledge Graph, and PostgreSQL for structured metadata and persistent state.
- **Orchestration:** Kubernetes (for agent lifecycle), Argo Workflows / Kubeflow Pipelines (for complex task orchestration).
- **Data Serialization:** Protocol Buffers (Protobuf) for rigorous schema enforcement and efficient data exchange.
- **Learning & Feedback:** Apache Flink/Spark Streaming (for real-time feedback processing), MLflow/DVC (for experiment tracking and data versioning), Ray RLlib (for policy optimization).

---

**Minute Aspects & Technical Dissection:**

### 1. Goal Decomposition & HTN Planning Engine Sub-system

This sub-system is the strategic mind that translates vague user aspirations into concrete, executable multi-agent action sequences.

- **1.1. High-Level Goal Interpreter & LLM-Driven Planner:**

  - **Purpose:** Takes the high-level intent (`TASK_REQUEST`, `COMPLEX_INFORMATION_QUERY`) and extracted entities from Nexus, and

interprets them into a structured, initial conceptual plan. This leverages a powerful LLM to handle the nuances of natural language goals.

- ○ **Technical Dissection:**
  - ■ **LLM Model:** A large, instruction-tuned LLM (e.g., Llama 3-70B, or a commercial API like GPT-4o if used via a secure, rate-limited adapter). This model is hosted on a dedicated inference cluster (e.g., leveraging NVIDIA TensorRT for optimization) and served by a high-throughput engine (e.g., vLLM).
  - ■ **System Prompt for Planning:** An incredibly detailed and dynamic system prompt is constructed. This prompt includes:
    - ■ **Role Definition:** "You are the central planning intelligence of Script Kitty. Your task is to break down complex user goals into a sequence of actionable steps, assigning them to the most appropriate internal agents (Armory, Skills, Guardian)."
    - ■ **Goal Context:** The user's original query, the intent identified by Nexus, and the `session_context` from Nexus.
    - ■ **Available Agent Capabilities:** A concise, semantically tagged summary of the *current* capabilities of Armory (tools, data sources), Skills (ML models, execution environments), and Guardian (policy checks, evaluation). This summary is dynamically retrieved from the Global Knowledge Graph (GKGS) and potentially updated by Foundry.
    - ■ **Planning Constraints/Policies:** High-level system constraints (e.g., "prioritize open-source tools," "minimize compute cost," "always ensure Guardian approval for external actions").
    - ■ **Desired Output Format:** A strict JSON schema (`PlanSchema` defined in Protobuf) outlining: `plan_id`, `goal_description`, `status`, `steps` (an array of `PlanStep` objects). Each `PlanStep` includes: `step_id`, `description`, `agent_assigned` (e.g., "Armory", "Skills"), `action_type` (e.g., "research", "execute_code", "evaluate"), `inputs` (parameters for the action), `expected_outputs`, `dependencies` (on other `step_id`s), and `max_retries`.
  - ■ **Reasoning Chain Generation (CoT/ToT):** The LLM is often prompted to generate a chain-of-thought or tree-of-thought internal reasoning process before outputting the final plan, increasing transparency and improving plan quality.
  - ■ **Output Validation & Correction:** The generated JSON plan is immediately validated against the `PlanSchema`. If invalid, the LLM is prompted for self-correction with specific error messages, or a symbolic planner (if used) takes over.

- **Fine-tuning & Learning:** The LLM's planning capability is continuously fine-tuned on pairs of user goals and human-approved, optimal execution plans, leveraging data from Guardian's feedback loop.
- **1.2. Hierarchical Task Network (HTN) & Symbolic Planning Fallback:**

  - **Purpose:** For well-defined, critical, or highly constrained tasks, or as a robust fallback for LLM planning failures, a symbolic HTN planner ensures formally correct and verifiable plan generation.
  - **Technical Dissection:**
    - **HTN Domain Definition:** A formal representation of Script Kitty's capabilities and action space is defined using an HTN-like language (e.g., a custom Python DSL or a subset of PDDL). This includes:
      - **Tasks:** Abstract goals (e.g., `research_topic`, `solve_programming_problem`).
      - **Methods:** Ways to achieve tasks (e.g., `research_topic` can be `web_scrape_and_summarize` via Armory, or `query_knowledge_graph` via GKGS).
      - **Operators:** Primitive, executable actions (e.g., `run_python_script`, `call_external_api`).
      - **Preconditions & Effects:** Formal logical statements describing what must be true before an action/method can start, and what changes occur after it completes.
    - **HTN Planner Implementation:** A custom Python implementation of an HTN planner or an off-the-shelf planner library (e.g., `py_pddl_planner`, if adapted) that can search the HTN domain for a valid plan.
    - **Plan Validation & Synthesis:** The symbolic planner generates a sequence of primitive actions. This plan is then translated into the `PlanSchema` structure, ensuring compatibility with the rest of Core's orchestration.
    - **Synergy with LLM:** The LLM can be used to interpret the user's initial goal into a symbolic state for the HTN planner, or the HTN planner can validate the LLM-generated plan for logical consistency.

## 2. Agent Orchestration & Lifecycle Management Sub-system

This sub-system is the operational core, dynamically summoning, monitoring, and managing the execution of tasks by specialized agents.

- **2.1. Task Scheduler & Dispatcher:**

  - **Purpose:** Based on the generated `PlanSchema`, this component schedules individual `PlanStep`s, ensuring dependencies are met and allocating tasks to the appropriate backend agents (Armory, Skills, Guardian).

- ○ **Technical Dissection:**
  - ■ **Directed Acyclic Graph (DAG) Execution:** The `PlanSchema` (with its `dependencies` field) is treated as a DAG. The scheduler identifies ready-to-execute steps (those with no unfulfilled dependencies).
  - ■ **Priority Queues:** Tasks are placed into priority queues based on urgency, resource requirements, or human override signals.
  - ■ **Agent Capability Matching:** Before dispatching a task, the scheduler performs a lookup in the GKGS to confirm the assigned agent's current availability and capability to execute the specific `action_type` defined in the `PlanStep`.
  - ■ **Message Generation:** For each ready `PlanStep`, a structured `TaskRequestSchema` (Protobuf) is generated, containing `task_id`, `step_id`, `agent_target`, `action_type`, `inputs` (extracted from the `PlanStep`), `callback_address` (for results), and `timeout`.
  - ■ **Communication:** Dispatches `TaskRequestSchema` messages to the appropriate agent's gRPC endpoint (e.g., Armory's `execute_research_task`, Skills' `run_code_execution`).
- ● **2.2. Agent Lifecycle & Resource Manager:**

  - ○ **Purpose:** Dynamically provisions or activates specialized agent instances within the Kubernetes cluster, monitors their health and resource consumption, and manages their graceful termination.
  - ○ **Technical Dissection:**
    - ■ **Kubernetes API Client:** Core leverages the official Kubernetes Python/Go client library to interact directly with the Kubernetes API server.
    - ■ **Dynamic Pod/Deployment Creation:** When a new `action_type` requires a specific, potentially ephemeral agent (e.g., a GPU-accelerated model for a complex vision task), Core issues API calls to create Kubernetes `Pod`s or `Deployment`s for that agent. The agent image and resource requirements (CPU, memory, GPU, custom volumes) are looked up from the GKGS's Agent Registry.
    - ■ **Resource Quotas & Limits:** Enforces resource quotas (at the namespace level) and pod-specific limits (CPU, memory) to prevent resource starvation or runaway processes.
    - ■ **Health Checks & Readiness Probes:** Continuously monitors the health and readiness of agent pods via Kubernetes `livenessProbe` and `readinessProbe` definitions. If an agent becomes unhealthy, Core can retry the task, reassign it, or flag it for human intervention.
    - ■ **Autoscaling Integration:** While Foundry handles cluster-wide autoscaling, Core can trigger specific horizontal pod autoscaling (HPA) events for its agents based on internal load or task queue depths.

- - ■ **Garbage Collection/Termination:** Automatically terminates agent pods/deployments once their assigned tasks are completed or after a predefined inactivity period, freeing up resources.
  - **2.3. Progress Tracking & State Update:**

    - ○ **Purpose:** Monitors the progress of ongoing tasks executed by agents, updates the global plan state, and relays progress updates to Nexus.
    - ○ **Technical Dissection:**
      - ■ **Internal State Machine:** Core maintains an internal state machine for each active `plan_id` and `step_id`, tracking its `status` (e.g., `PENDING`, `DISPATCHED`, `RUNNING`, `COMPLETED`, `FAILED`, `CANCELLED`).
      - ■ **Asynchronous Callbacks:** Agents report their status and results back to Core via dedicated gRPC endpoints (`report_task_update`, `report_task_result`).
      - ■ **Timeout & Retries:** Each `PlanStep` has a defined `timeout`. If a task exceeds this, Core can automatically trigger a retry (up to `max_retries`) or mark the step as `FAILED`.
      - ■ **Error Handling & Fallback:** Comprehensive error handling for agent failures. Core determines if a task can be retried, requires a different approach (triggering replanning by the LLM), or needs human intervention (notifying Guardian).
      - ■ **Nexus Update Dispatch:** Periodically (e.g., every 5 seconds for long tasks) or upon significant state changes, Core sends `TaskUpdateSchema` messages to Nexus via the internal communication hub, containing `plan_id`, `step_id`, `status`, and a `message` for the user (e.g., "Armory is currently researching...")
      - ■ **Transactionality (Lightweight):** For critical state changes (e.g., marking a step complete), Core might use atomic database operations or distributed transactions to ensure data consistency.

## 3. Global Knowledge Graph System (GKGS) & Context Store Sub-system

The GKGS is the collective, evolving memory and "common sense" of Script Kitty, a central source of truth for all components.

- **3.1. Knowledge Graph (KG) Database:**

  - ○ **Purpose:** The central, authoritative repository for all Script Kitty's acquired knowledge, facts, rules, learned heuristics, tool manifests, agent capabilities, and ongoing task contexts.
  - ○ **Technical Dissection:**

- **Graph Database:** Neo4j Community Edition or ArangoDB (multi-model) is the preferred choice due to their native support for relationships and efficient traversal.
- **Schema Definition:** A rigorous ontology/schema defines node types (e.g., `User`, `Concept`, `Tool`, `Agent`, `Task`, `Policy`) and relationship types (e.g., `HAS_CAPABILITY`, `USES_TOOL`, `IS_A`, `DEPENDS_ON`, `DEFINES_POLICY`, `HAS_PROPERTY`).
- **Data Model:** Knowledge is stored as triples (`subject`, `predicate`, `object`) or property graphs (`nodes` with `properties`, `edges` with `properties`). This allows for highly flexible and expressive knowledge representation.
- **Persistent Storage:** Data is persisted to disk with journaling and replication for high availability.
- **Query Language:** Cypher (for Neo4j) or AQL (for ArangoDB) is used for complex graph traversals and pattern matching to retrieve relevant knowledge.
- **Real-time Updates:** Exposed gRPC endpoints (`add_knowledge`, `update_knowledge`, `delete_knowledge`) allow Core and other agents (especially Armory and Guardian) to dynamically update the KG.

- **3.2. Vector Store & Semantic Search:**

  - **Purpose:** Enables semantic search and Retrieval-Augmented Generation (RAG) over the knowledge graph, allowing LLMs and other agents to retrieve contextually relevant information using natural language queries.
  - **Technical Dissection:**
    - **Embedding Models:** Text chunks (e.g., descriptions of tools, summaries of concepts, policy rules) from the KG are converted into high-dimensional vector embeddings using state-of-the-art embedding models (e.g., `all-MiniLM-L6-v2` or `OpenAI Embeddings` if using an API).
    - **Vector Database:** Weaviate or ChromaDB are integrated to store these vector embeddings along with references back to their original nodes in the Neo4j KG.
    - **Indexing:** Efficient indexing (e.g., HNSW for approximate nearest neighbor search) ensures low-latency retrieval.
    - **RAG Pipeline:** When an LLM in Core needs external knowledge for planning or reasoning, it generates a query embedding. The vector store performs a semantic search to find the most relevant chunks. These chunks are then inserted into the LLM's prompt as context.
    - **Knowledge Graph Embedding (KGE) Learning (Future):** For advanced reasoning, KGE techniques (e.g., TransE, ComplEx via libraries like `pykeen`) can learn embeddings of entities and relations *within* the graph itself, enabling inferential reasoning and link prediction.

- **3.3. Knowledge Ingestion & Validation:**

  - **Purpose:** Manages the process of adding new knowledge to the GKGS, whether from Armory's research, human feedback, or system-generated insights.
  - **Technical Dissection:**
    - **Schema Enforcement:** New knowledge entries are rigorously validated against the KG's schema to ensure structural integrity and prevent inconsistencies.
    - **Conflict Resolution:** Mechanisms for handling conflicting information, potentially involving confidence scores, timestamps, or human arbitration.
    - **Version Control (Semantic):** While not full Git for graphs, a semantic versioning system for knowledge chunks can be implemented, allowing rollbacks or tracking of changes.
    - **Automated Graph Completion (Optional):** Techniques like rule-based reasoning or even small, specialized LLMs can infer new relationships or facts from existing knowledge, enriching the graph.

## 4. Internal Communication Hub & Message Broker Sub-system

This is the nervous system of Script Kitty, ensuring efficient and reliable communication between Core and all specialized backend agents.

- **4.1. gRPC Server & Client (Core-Agent Communication):**

  - **Purpose:** Provides a high-performance, strongly-typed, and bi-directional communication channel for direct request-response interactions between Core and its agents (Armory, Skills, Guardian).
  - **Technical Dissection:**
    - **gRPC Server:** Core runs a gRPC server, exposing a set of RPC methods (defined in Protobuf service definitions) for agents to call (e.g., `report_task_update`, `report_task_result`).
    - **gRPC Clients:** Core itself acts as a gRPC client, initiating calls to agents (e.g., `Armory.execute_research_task`, `Skills.run_code`).
    - **Protocol Buffers (Protobuf):** The cornerstone for data exchange. Every message and service contract is precisely defined in `.proto` files. This ensures:
      - **Strong Typing:** Prevents common integration errors by enforcing data types.
      - **Backward/Forward Compatibility:** Allows independent deployment of services with different versions of message schemas (within compatible changes).
      - **Efficiency:** Highly optimized binary serialization/deserialization, minimizing payload size and processing time.

- ■ **Language Agnostic:** Generated code for Go, Python, etc., ensures seamless inter-language communication.
  - ■ **Streaming RPCs:** Utilized for scenarios like continuous progress updates from a long-running Skill task to Core, or Core sending a stream of micro-instructions to an agent.
  - ■ **Error Handling:** Standard gRPC error codes are used, providing clear communication of failures between services.
- ● **4.2. Asynchronous Event Bus (Kafka/Pulsar):**

  - ○ **Purpose:** Provides a decoupled, scalable, and fault-tolerant mechanism for asynchronous event streaming, crucial for audit logs, performance metrics, raw experience data, and background updates.
  - ○ **Technical Dissection:**
    - ■ **Apache Kafka / Apache Pulsar Cluster:** A robust distributed streaming platform is deployed.
    - ■ **Topics:** Various Kafka topics are defined for different types of events (e.g., `script_kitty_task_events`, `script_kitty_metrics`, `script_kitty_audit_logs`, `script_kitty_raw_experience_data`).
    - ■ **Producers:** All agents and Core components act as producers, asynchronously publishing events to relevant topics.
    - ■ **Consumers:** Core, Guardian (for evaluation), and Foundry (for MLOps) act as consumers, subscribing to relevant topics to process event streams.
    - ■ **Durability & Replication:** Messages are durably persisted and replicated across brokers, ensuring no data loss even in case of node failures.
    - ■ **Scalability:** The distributed nature of Kafka/Pulsar allows for handling massive volumes of events and scaling consumption dynamically.
    - ■ **Backpressure Handling:** Consumers can process events at their own pace without overwhelming producers, providing system stability.

## 5. Feedback & Learning Integration Sub-system

This is the engine of Script Kitty's continuous improvement, absorbing performance data, user feedback, and success/failure signals to refine its own strategic capabilities.

- ● **5.1. Performance Monitoring & Anomaly Detection (Core Integration):**

  - ○ **Purpose:** Ingests real-time metrics and performance data from Foundry's observability stack, enabling Core to react to degradation or anomalies in agent performance.
  - ○ **Technical Dissection:**

- **Metrics Consumption:** Core subscribes to `script_kitty_metrics` Kafka topics, consuming aggregated metrics like agent latency, error rates, resource utilization, and task success rates.
- **Internal Anomaly Detection:** Core can run lightweight anomaly detection algorithms (e.g., moving averages, simple statistical models) on these streams to identify immediate issues with agents or specific task types.
- **Adaptive Behavior:** If an anomaly is detected (e.g., a specific Skill agent consistently failing), Core can adapt its planning (e.g., avoid that agent for a short period, trigger an auto-restart via Agent Lifecycle Manager, or escalate to Guardian for deeper analysis/retraining).

- **5.2. Learning Data Ingestion & Experience Replay Buffer:**

  - **Purpose:** Collects structured data on task execution outcomes (successes, failures, human corrections), forming the basis for improving planning heuristics and model training.
  - **Technical Dissection:**
    - **Structured Experience Data:** When a task completes (successfully or not), Core constructs a `TaskExperienceSchema` (Protobuf) containing: `plan_id`, `original_goal`, `executed_plan` (the sequence of steps), `final_result`, `human_feedback` (if any), `metrics` (e.g., execution time, resources used), and `success_boolean`.
    - **Kafka Stream:** These `TaskExperienceSchema` messages are published to a dedicated `script_kitty_experience_data` Kafka topic.
    - **Experience Replay Buffer (Persistent Storage):** A consumer for this topic persists the experience data to a long-term storage solution (e.g., object storage like MinIO, or a data lake like Parquet files on S3) and also loads it into an in-memory or Redis-backed "replay buffer" for immediate use in policy learning.

- **5.3. Planning Heuristic Refinement & Policy Learning:**

  - **Purpose:** Utilizes the collected experience data to continuously improve Core's planning strategies, making it more efficient, reliable, and intelligent over time.
  - **Technical Dissection:**
    - **Reinforcement Learning (RL) Pipeline:** For complex, sequential decision-making in planning, Core can employ an RL-based approach.
      - **State:** The current `session_context`, `goal`, and `available_agents_capabilities`.
      - **Actions:** The possible `PlanStep`s to choose from and the agent to assign.
      - **Reward Function:** Defined based on task success, efficiency (time, cost), adherence to policies (penalties for violations), and positive human feedback.

- - - **RL Algorithm:** An off-policy algorithm like PPO or SAC (using Ray RLlib) can be used to train a policy network that learns optimal planning strategies from the experience data.
  - **Imitation Learning/Supervised Learning:** For more direct refinement, Core can use supervised learning on human-corrected plans, or imitation learning from expert demonstrations (successful plans).
  - **Feature Engineering:** Features for the learning model include attributes of the goal, the agents, and the context, derived from the GKGS.
  - **Model Training:** The RL policy or supervised model for planning refinement is periodically retrained (triggered by Guardian or scheduled by Foundry) using the latest experience data.
  - **MLflow Integration:** MLflow tracks experiments, logging hyperparameters, model metrics, and artifacts for different planning policy versions, enabling reproducible research and model selection.
  - **Policy Deployment:** New, improved planning policies are deployed via Foundry's CI/CD pipelines as part of Core's updated container images, ensuring seamless integration.
- **5.4. Knowledge Graph Update Triggers:**

  - **Purpose:** Based on new insights from learning or direct directives from Guardian, Core triggers updates to the GKGS, ensuring the system's "common sense" evolves.
  - **Technical Dissection:**
    - **Automated Knowledge Extraction:** From successful task experiences, Core can automatically extract new facts or relationships (e.g., "Tool X was successfully used for Task Y," "Concept Z is related to Concept A"). This extraction can involve small, specialized LLMs or rule-based parsers.
    - **Guardian Directives:** Guardian, after detecting a policy violation or identifying a new best practice, can send specific instructions to Core to update policies or capabilities in the GKGS.
    - **GKGS API Calls:** Core uses the GKGS's gRPC API (`add_knowledge`, `update_knowledge`) to programmatically inject these new insights, ensuring the GKGS remains the single source of truth.

---

**Operational Excellence & Scalability (Core-Specific):**

- **High Availability:** Core is designed as a distributed, stateless (where possible) microservice system. Redundant Core instances are deployed across multiple Kubernetes nodes and availability zones.
- **Load Balancing:** Service meshes (Istio/Linkerd) provide intelligent load balancing for gRPC traffic to Core's services, ensuring requests are distributed efficiently.

- **Resilience:** Circuit breakers and bulkheads are implemented in gRPC client calls from Core to agents, preventing cascading failures if a downstream agent becomes unresponsive.
- **Observability:** Comprehensive logging, metrics, and tracing (via OpenTelemetry) provide deep visibility into Core's planning process, agent dispatch, and overall system health. Dashboards in Grafana allow real-time monitoring of planning success rates, agent utilization, and communication latencies.
- **Dynamic Configuration:** Core's operational parameters (e.g., retry counts, timeouts, LLM temperature for planning) are managed via a centralized configuration service (e.g., Kubernetes ConfigMaps or a dedicated service like Consul/etcd), allowing for dynamic updates without service restarts.
- **Fault Tolerance:** The use of Kafka/Pulsar for event streaming provides inherent fault tolerance for asynchronous communication, as messages are durably stored.
- **Horizontal Scalability:** Most Core components (Task Scheduler, Agent Lifecycle Manager) are designed to be horizontally scalable, adding more instances as the number of concurrent plans increases. The LLM inference service for planning scales independently.
- **Cost Optimization:** The Agent Lifecycle Manager actively monitors resource usage and dynamically spins down idle agent pods, minimizing compute costs. The planning LLM instance itself can be scaled based on demand, potentially using spot instances for non-critical planning tasks.

# Core interactions

**Core Communication Paradigms and Infrastructure:**

At the heart of Script Kitty's inter-component communication lies a meticulously designed infrastructure built for high performance, reliability, and security.

1. **gRPC (Remote Procedure Call):**

   ○ **Purpose:** Primary mechanism for synchronous, high-throughput, point-to-point communication where a request expects an immediate response. Ideal for task initiation, status inquiries, and result retrieval.
   ○ **Technical Dissection:**
      ■ **Protocol Buffers (Protobuf):** All gRPC messages are defined using Protobuf schemas (e.g., `ScriptKittyMessageSchema`, `TaskRequestSchema`, `ToolManifestSchema`). This ensures strict type-checking, efficient binary serialization/deserialization, and language-agnostic interface definitions. Each field, its type, and its required/optional status are precisely defined, eliminating ambiguity.
      ■ **HTTP/2:** gRPC leverages HTTP/2 for multiplexing (multiple concurrent requests over a single connection), header compression, and stream-based communication, significantly reducing overhead compared to traditional REST over HTTP/1.1.
      ■ **Server-Side Streaming:** Used for scenarios where a server (e.g., Core) needs to send multiple messages back to a client (e.g., Nexus sending periodic task updates).
      ■ **Client-Side Streaming & Bidirectional Streaming:** Potentially for complex, interactive negotiations between agents (e.g., Core providing a stream of sub-tasks to Skills, which in turn streams back progress updates).
      ■ **Load Balancing & Service Discovery:** Integrated with Foundry's Kubernetes service mesh (Istio/Linkerd) for intelligent load balancing across multiple service instances, and Kube-DNS for efficient service discovery.
      ■ **Error Handling:** Standard gRPC status codes are used for synchronous error reporting (e.g., `UNAVAILABLE`, `DEADLINE_EXCEEDED`, `PERMISSION_DENIED`). Retries with exponential backoff are implemented at the client level.
      ■ **Security:** Enforced via mTLS (Mutual TLS) provided by Foundry's service mesh, ensuring all inter-service gRPC communication is encrypted and mutually authenticated.
2. **Apache Kafka (Event Streaming Platform):**

- ○ **Purpose:** The backbone for asynchronous, decoupled, and highly scalable event-driven communication. Ideal for progress updates, telemetry, audit logs, raw experience data streams, and triggering downstream processes.
- ○ **Technical Dissection:**
  - ■ **Topics:** Messages are organized into topics (e.g., `sk_task_updates`, `sk_agent_metrics`, `sk_audit_logs`, `sk_raw_experience`). Each topic is partitioned across multiple brokers for scalability and fault tolerance.
  - ■ **Producers:** Script Kitty agents (e.g., Skills, Armory, Guardian) publish messages to relevant topics. Producers are designed to be idempotent and handle retries for transient network issues.
  - ■ **Consumers:** Other agents or Foundry components (e.g., Guardian for evaluation, Foundry for metrics ingestion) subscribe to topics and consume messages. Consumer groups ensure messages are processed exactly-once within a group.
  - ■ **Schema Registry (Confluent Schema Registry or compatible):** Ensures all Kafka messages adhere to predefined Protobuf schemas, providing strong data governance and preventing schema evolution issues.
  - ■ **Durability & Retention:** Messages are durably stored on disk for configurable retention periods, allowing consumers to replay events or catch up after downtime.
  - ■ **Fault Tolerance:** Kafka's distributed architecture (replication, leader/follower model) provides high availability even in the event of broker failures.
  - ■ **Security:** Kafka uses SASL (Simple Authentication and Security Layer) for authentication (e.g., Kerberos, SCRAM) and TLS for encryption of data in transit. Authorization is managed via ACLs (Access Control Lists) to restrict which agents can read/write to specific topics.
3. **RESTful APIs (Limited & External-Facing):**

- ○ **Purpose:** Primarily used for specific interactions with external services where gRPC is not supported, or for certain idempotent, read-heavy operations within Foundry.
- ○ **Technical Dissection:**
  - ■ **HTTP/S:** Standard HTTP/S protocol.
  - ■ **JSON/Protobuf:** Payloads are typically JSON, but can also be Protobuf for internal APIs.
  - ■ **Authentication/Authorization:** OAuth2, API Keys, or token-based authentication.
  - ■ **Rate Limiting:** Implemented to prevent abuse and manage load.
4. **Shared Global Knowledge Graph (GKGS) - Centralized Data Store:**

- ○ **Purpose:** While not a communication protocol, the GKGS (managed by Core and Foundry) serves as the central, authoritative source of shared system knowledge, context, and long-term memory. Agents interact with it through Core's API or direct read access (if authorized). It acts as a persistent, high-bandwidth "communication channel" for implicit knowledge transfer.
- ○ **Technical Dissection:** Graph database (Neo4j/ArangoDB) for structured knowledge, Vector Database (Weaviate/Chroma) for semantic search. Access is governed by fine-grained permissions.

---

### Inter-Component Interactions: Core and Others

**Script Kitty.Core** is the brain, the high-level orchestrator. Its interactions are central to all operations, acting as the primary hub for task decomposition, resource allocation, and progress monitoring.

---

### 1. Script Kitty.Core ↔ Script Kitty.Nexus: The Command & Response Loop

This is the primary user-facing communication channel, driven by the Nexus's role as the intermediary.

- ● **1.1. Nexus to Core: User Intent & Task Request Submission**

  - ○ **Purpose:** Nexus translates a user's natural language query into a structured `TaskRequestSchema` and sends it to Core for processing.
  - ○ **Flow & Data Schema:**
    - ■ **Nexus (Core NLU & Intent Routing):** Receives `ScriptKittyMessageSchema` (raw text + session context) from its User Input & Pre-processing sub-system.
    - ■ The Nexus LLM (Mistral-7B/Llama 3-8B) analyzes the `raw_text` and `session_context` to identify the `intent` (e.g., `TASK_REQUEST`, `COMPLEX_INFORMATION_QUERY`) and `extracted_entities` (e.g., "summarize," "research topic," "code for X").
    - ■ If the router determines the request requires Core's planning capabilities, it constructs a `TaskRequestSchema` (Protobuf) containing:
      - ■ `task_id` (UUID, globally unique)
      - ■ `session_id` (from Nexus session context)
      - ■ `user_id`
      - ■ `original_query` (raw text)
      - ■ `parsed_intent` (enum: `PLAN_TASK`, `RETRIEVE_INFO`, etc.)
      - ■ `parsed_entities` (map/list of key-value pairs)

- ■ `requester_agent_id` (Nexus)
- ■ `priority` (e.g., `HIGH`, `MEDIUM`)
- ■ `deadline` (if specified by user or inferred)
- ○ **Communication Protocol:** gRPC. Nexus's `Core Communication Gateway` (gRPC client) invokes a `ProcessTaskRequest` method on Core's `Internal Communication Hub` (gRPC server).
- ○ **Latency & Reliability:** Critical path for user experience. Nexus employs gRPC client-side retries with exponential backoff for transient Core unavailability. Core's gRPC server is highly available and fault-tolerant, running multiple replicas managed by Foundry.
- **1.2. Core to Nexus: Task Planning Updates (Asynchronous)**

  - ○ **Purpose:** Core provides Nexus with real-time updates on the progress of its task planning and sub-task delegation, maintaining the "group chat" illusion.
  - ○ **Flow & Data Schema:**
    - ■ **Core (Goal Decomposition & HTN Planning Engine):** After receiving a `TaskRequest`, Core's planning engine initiates its work. It generates intermediate states, identifies initial sub-goals, and assigns them.
    - ■ Core periodically (or upon significant state changes) publishes `TaskUpdateSchema` (Protobuf) messages to a Kafka topic (`sk_task_updates`). This schema includes:
      - ■ `task_id`
      - ■ `session_id`
      - ■ `update_type` (enum: `PLANNING_STARTED`, `SUBTASK_DELEGATED`, `PLANNING_PROGRESS`, `PLANNING_FAILED`)
      - ■ `current_status_message` (human-readable string)
      - ■ `progress_percentage` (optional)
      - ■ `sub_task_details` (e.g., `agent_id`, `sub_task_description`)
      - ■ `source_agent_id` (`Core`)
  - ○ **Communication Protocol:** Kafka. Core's `Internal Communication Hub` (Kafka producer) publishes to `sk_task_updates`.
  - ○ **Nexus (Intermediary & Internal Communication Proxy):** Nexus's `Core Communication Gateway` (Kafka consumer) subscribes to `sk_task_updates`.
  - ○ **Nexus (Response Synthesis & Persona Layer):** Upon receiving a `TaskUpdate`, Nexus's persona layer synthesizes a concise, persona-consistent message for the user, possibly prefixing with "Core is planning..." or "Core has

delegated a sub-task...". This update is then rendered and dispatched via the `User Output Renderer & Channel Dispatcher`.

- ○ **Asynchronicity:** Decouples Core's complex planning from Nexus's immediate UI responsiveness, allowing Core to take its time without blocking the user interface.
- **1.3. Core to Nexus: Final Task Result Delivery**

  - ○ **Purpose:** Core delivers the final result of a completed task back to Nexus for presentation to the user.
  - ○ **Flow & Data Schema:**
    - ■ **Core (Agent Orchestration & Lifecycle Management):** Once all sub-tasks orchestrated by Core are completed and the final result is aggregated, Core prepares the `TaskResultSchema` (Protobuf). This schema contains:
      - ■ `task_id`
      - ■ `session_id`
      - ■ `result_status` (enum: `SUCCESS`, `FAILURE`, `PARTIAL_SUCCESS`)
      - ■ `final_output` (structured data, e.g., text summary, code snippet, image URL, nested JSON results)
      - ■ `execution_summary` (optional, e.g., steps taken, agents involved, duration)
      - ■ `source_agent_id` (`Core`)
    - ■ **Core (Internal Communication Hub):** Publishes `TaskResultSchema` to the `sk_task_results` Kafka topic.
  - ○ **Communication Protocol:** Kafka.
  - ○ **Nexus (Core Communication Gateway):** Consumes from `sk_task_results`.
  - ○ **Nexus (Response Synthesis & Persona Layer):** Synthesizes the `final_output` into a coherent, persona-driven conversational response, attributing it appropriately (e.g., "The Nexus has completed your request:"). This is then dispatched to the user.
  - ○ **Error Handling:** If `result_status` is `FAILURE`, Nexus's persona layer synthesizes an appropriate error message, perhaps suggesting rephrasing the query or indicating a system issue, potentially logging the failure for Guardian's review.

---

## 2. Script Kitty.Core ↔ Script Kitty.Armory: Resource Acquisition & Tooling

Core relies on Armory for access to external information and dynamic tool management.

- **2.1. Core to Armory: Tool Discovery & Request for Resource/Research**

- **Purpose:** Core, during its planning phase, identifies the need for a specific tool or external research (e.g., "find current stock prices," "scrape data from URL," "provision GPU for model"). It queries Armory for available capabilities.
- **Flow & Data Schema:**
    - **Core (Goal Decomposition & HTN Planning Engine):** Based on its plan and sub-goals, Core determines a specific external action is needed. It queries the Global Knowledge Graph (GKGS) for known tool capabilities or sends a semantic query to Armory.
    - **Core (Internal Communication Hub):** Invokes gRPC calls on Armory's API:
        - `DiscoverTools(semantic_query: string, required_capabilities: list<string>)`: Core requests tools that can perform a certain type of action (e.g., `semantic_query="web scraping"`, `required_capabilities=["extract_text", "handle_javascript"]`).
        - `RequestResource(resource_type: enum, config: map<string,string>)`: Requests provisioning of compute or data resources (e.g., `resource_type="GPU_VM"`, `config={"gpu_type": "A100", "duration_hours": 2}`).
        - `RequestResearch(research_query: string, target_domains: list<string>)`: Requests Armory to perform web scraping and data digestion for a specific query.
    - **Armory (Dynamic Tool Registry & Discovery):** Armory's `Tool Registry` service receives `DiscoverTools` requests. It performs a semantic search over its `ToolManifestSchema` database (PostgreSQL/Elasticsearch), matching `semantic_query` and `required_capabilities` against registered tool descriptions. It returns a list of matching `ToolManifestSchema` objects (containing tool name, ID, input/output schemas, required env).
    - **Armory (Intelligent Web Scraping & Data Digestion):** For `RequestResearch`, Armory's `Scraping/Digestion` sub-system initiates the research process.
    - **Armory (Resource Provisioning via IaC Templates):** For `RequestResource`, Armory's `IaC Orchestrator` component processes the request.
- **Communication Protocol:** gRPC for requests.
- **Response from Armory:** Armory's gRPC server returns `ToolDiscoveryResult` (list of `ToolManifestSchema`), `ResourceProvisioningStatus` (success/failure, resource ID, access details), or `ResearchTaskID`.

- **2.2. Armory to Core: Tool/Resource Provisioning Status & Research Results**

  - **Purpose:** Armory reports back to Core on the status of resource provisioning or delivers the results of a research task.
  - **Flow & Data Schema:**
    - **Armory (Resource Provisioning/Web Scraping):** Asynchronously processes requests. Once a resource is provisioned (or fails), or research is completed, Armory publishes an update.
    - **Armory (Internal Communication Hub):** Publishes `ArmoryUpdateSchema` (Protobuf) messages to a Kafka topic (`sk_armory_updates`). This schema includes:
      - `task_id` (correlates with Core's original request)
      - `update_type` (enum: `RESOURCE_PROVISIONED`, `RESOURCE_FAILURE`, `RESEARCH_COMPLETED`, `RESEARCH_FAILURE`, `TOOL_WRAPPER_GENERATED`)
      - `details` (e.g., `resource_id`, `access_credentials`, `research_summary`, `link_to_digested_data_in_data_lake`, `tool_manifest_id`)
      - `source_agent_id` (`Armory`)
  - **Communication Protocol:** Kafka (asynchronous updates).
  - **Core (Agent Orchestration & Lifecycle Management):** Subscribes to `sk_armory_updates`. Upon receiving an update, Core's `Orchestrator` updates the plan's state in its internal memory and the GKGS (e.g., marking a resource as available, integrating research findings).
  - **GKGS Integration:** The `research_summary` and `digested_data_link` are integrated into the GKGS by Core's `GKGS & Context Store` component, making the newly acquired knowledge universally accessible for future planning and RAG.
- **2.3. Core to Armory: Tool Wrapper Generation (Human-Verified)**

  - **Purpose:** When Core determines a needed tool lacks a ready-to-use wrapper, or a custom one is needed, it requests Armory to generate code for it.
  - **Flow & Data Schema:**
    - **Core (HTN Planning Engine):** Identifies a missing tool interface for a specific capability. Sends a `GenerateToolWrapperRequestSchema` via gRPC to Armory, including:
      - `tool_manifest_id` (or `tool_description_text`)
      - `target_language` (e.g., `PYTHON`)
      - `desired_api_signature` (e.g., function name, parameters, return type)

- - - **example_usage** (optional)
  - **Armory (Templated API/CLI Wrapper Generation):** Uses its code-focused LLM and templating engine to generate the wrapper code.
  - **Armory (Human Review Interface Integration):** Presents the generated code for human review and approval.
  - **Armory to Core (via Kafka):** Publishes an `ArmoryUpdateSchema` (type `TOOL_WRAPPER_GENERATED`) with the `tool_manifest_id` and confirmation that the wrapper is ready for deployment (post-human review). Foundry's CI/CD would then deploy this new wrapper.

---

**3. Script Kitty.Core ↔ Script Kitty.Skills: Task Execution & Specialized AI**

Core delegates the actual execution of computational tasks and AI model inference to Skills agents.

- **3.1. Core to Skills: Task Execution Delegation**

  - **Purpose:** Core breaks down its high-level plan into concrete, executable sub-tasks and delegates them to specific Skills agents for execution.
  - **Flow & Data Schema:**
    - **Core (Agent Orchestration & Lifecycle Management):** Based on the HTN plan, Core identifies `action_nodes` that require execution by Skills. It looks up the `ToolManifestSchema` (from GKGS/Armory) for the required Skill (e.g., `CodeExecutionSkill`, `ImageRecognitionSkill`, `DataAnalysisSkill`).
    - Core constructs a `SkillExecutionRequestSchema` (Protobuf) containing:
      - `task_id` (parent task ID from Nexus)
      - `sub_task_id` (unique ID for this specific Skills execution)
      - `skill_id` (ID of the specific Skill agent/model to invoke)
      - `input_parameters` (structured data matching the Skill's manifest, e.g., `{"code_to_execute": "...", "language": "python"}`, or `{"image_url": "...", "model_version": "v2"}`)
      - `output_schema_expected` (optional, for validation)
      - `resource_requirements` (e.g., `{"gpu_required": true, "memory_gb": 16}`)
      - `execution_timeout_seconds`
      - `source_agent_id` (`Core`)

- ■ **Core (Internal Communication Hub):** Invokes a gRPC call on the target Skills agent's `Skills API Gateway` (e.g., `ExecuteSkillTask`).
  - ○ **Communication Protocol:** gRPC.
  - ○ **Foundry's Role:** Foundry's Kubernetes-Native Compute & Orchestration, through the custom `SkillModule` CRD, would ensure the correct Skills pod is running or spun up with the necessary resources (e.g., a GPU-enabled pod for an LLM-based code generation skill).
- ● **3.2. Skills to Core: Task Execution Updates (Asynchronous)**

  - ○ **Purpose:** Skills reports granular progress and intermediate results during long-running tasks, allowing Core to monitor execution and adapt its plan if necessary.
  - ○ **Flow & Data Schema:**
    - ■ **Skills (Sandboxed Code Execution/Specialized AI Models):** During execution, Skills agents publish `SkillUpdateSchema` (Protobuf) messages to a Kafka topic (`sk_skill_updates`). This schema includes:
      - ■ `task_id`
      - ■ `sub_task_id`
      - ■ `update_type` (enum: `STARTED`, `PROGRESS`, `LOG_MESSAGE`, `INTERMEDIATE_RESULT`, `ERROR`)
      - ■ `status_message`
      - ■ `progress_percentage`
      - ■ `payload` (e.g., partial data, log line, error trace)
      - ■ `source_agent_id` (specific Skill instance ID)
  - ○ **Communication Protocol:** Kafka. Skills' internal `Communication Module` (Kafka producer) publishes to `sk_skill_updates`.
  - ○ **Core (Agent Orchestration & Lifecycle Management):** Core's `Orchestrator` (Kafka consumer) subscribes to `sk_skill_updates`. It updates the state of the relevant `sub_task_id` in its internal `TaskState` store.
  - ○ **Core's Decision Making:** Core uses these updates for adaptive planning:
    - ■ If `progress_percentage` stalls, it might investigate or retry.
    - ■ If `update_type` is `ERROR`, Core can:
      - ■ Attempt a different strategy (e.g., try another Skill agent, reformulate the input).
      - ■ Delegate a debugging sub-task to a `DebuggingSkill`.
      - ■ Report back to Nexus about a failure (which then goes to the user).
      - ■ Log a detailed failure for Guardian's `Automated Evaluation`.
- ● **3.3. Skills to Core: Final Task Execution Result**

  - ○ **Purpose:** Skills delivers the definitive outcome of its assigned sub-task.

- ○ **Flow & Data Schema:**
  - ■ **Skills (Sandboxed Code Execution/Specialized AI Models):** Upon completion or final failure, Skills constructs `SkillResultSchema` (Protobuf) and publishes it to `sk_skill_results` Kafka topic.
    - ■ `task_id`
    - ■ `sub_task_id`
    - ■ `result_status` (enum: `SUCCESS`, `FAILED`, `TIMED_OUT`)
    - ■ `output_data` (structured result, e.g., code output, analysis report, image analysis results)
    - ■ `error_details` (if failed)
    - ■ `execution_metrics` (e.g., CPU time, memory used, model inference time)
    - ■ `source_agent_id` (specific Skill instance ID)
- ○ **Communication Protocol:** Kafka.
- ○ **Core (Agent Orchestration & Lifecycle Management):** Consumes from `sk_skill_results`.
- ○ **Core's Plan Execution:** Core updates the `sub_task_id` as `COMPLETED`. If successful, it integrates `output_data` into its next planning step or aggregates it towards the final `TaskResultSchema` to be sent to Nexus. If failed, it triggers error handling within its planning engine.
- ○ **Guardian Integration (via Foundry's Observability):** `execution_metrics` and `error_details` are automatically collected by Foundry's observability stack and consumed by Guardian for performance monitoring, anomaly detection, and automated evaluation/training triggers.

---

**4. Script Kitty.Core ↔ Script Kitty.Guardian: Policy, Safety & Learning Feedback**

Guardian acts as Script Kitty's conscience and continuous improvement engine, directly influencing Core's actions and learning processes.

- ● **4.1. Core to Guardian: Action Approval Request**

  - ○ **Purpose:** For critical or potentially sensitive actions identified in Core's plan, Core *must* request approval from Guardian's Policy Enforcement Engine before execution. This is a crucial safety and alignment mechanism, potentially involving human-in-the-loop.
  - ○ **Flow & Data Schema:**
    - ■ **Core (Agent Orchestration & Lifecycle Management):** During plan generation, Core's `HTN Planning Engine` identifies sensitive `action_nodes` (e.g., "accessing external API with write privileges,"

"performing an action affecting real-world systems," "generating code that might interact with OS," "accessing PII"). These actions are flagged.

- Core constructs an `ActionApprovalRequestSchema` (Protobuf) and sends it via gRPC to Guardian.
  - `task_id`
  - `action_id` (unique ID for this specific action)
  - `proposed_action_description` (natural language summary of the action)
  - `action_parameters` (structured input to the action)
  - `risk_assessment_by_core` (Core's own confidence/risk assessment, optional)
  - `required_policy_tags` (e.g., `DATA_PRIVACY`, `EXTERNAL_ACCESS`, `SAFETY_CRITICAL`)
  - `source_agent_id` (`Core`)

- **Communication Protocol:** gRPC. Core's `Internal Communication Hub` invokes `RequestActionApproval` on Guardian's `Policy Enforcement Engine` (gRPC server).

- **Guardian (Policy Enforcement Engine):**
  - Receives the request.
  - OPA (Open Policy Agent) evaluates the `proposed_action_description` and `action_parameters` against predefined Rego policies.
  - Its ethical LLM might provide an initial ethical assessment.
  - If high risk or policy violation is detected, Guardian might trigger a human review via its `Human Oversight & Feedback Interface`.

- **Response from Guardian:** Guardian returns `ActionApprovalResponseSchema` (Protobuf) via gRPC synchronous response:
  - `action_id`
  - `approval_status` (enum: `APPROVED`, `DENIED`, `PENDING_HUMAN_REVIEW`)
  - `reason` (if denied)
  - `policy_violations` (list of violated policies)

- **Core's Reaction:** If `APPROVED`, Core proceeds with plan execution. If `DENIED`, Core's planning engine initiates a fallback (e.g., tries an alternative path, reframes the action, reports to Nexus). If `PENDING_HUMAN_REVIEW`, Core pauses that branch of the plan and waits for a later update from Guardian (via Kafka).

- **4.2. Guardian to Core: Policy Updates & Feedback for Learning**

- ○ **Purpose:** Guardian informs Core about policy changes, provides direct feedback on Core's planning performance, and triggers updates to Core's planning heuristics.
- ○ **Flow & Data Schema:**
  - ■ **Guardian (Automated Evaluation & Training Trigger):**
    - ■ **Policy Updates:** When human operators update policies via Guardian's UI, Guardian pushes new `PolicyUpdateSchema` (Protobuf) messages to a Kafka topic (`sk_policy_updates`).
    - ■ **Performance Feedback:** Based on its `Real-time Performance Monitoring` (metrics, logs, traces from Foundry) and `Automated Evaluation`, Guardian identifies instances of Core's successful or unsuccessful planning, efficient or inefficient resource allocation. It synthesizes this into `FeedbackSchema` (Protobuf) messages for Core's learning.
    - ■ **GKGS Updates:** Guardian can directly instruct Core to update the GKGS based on newly learned facts or refined heuristics (e.g., `GKGSUpdateSchema` via gRPC).
  - ■ **Communication Protocol:** Kafka for asynchronous updates (`sk_policy_updates`, `sk_feedback_for_core`). gRPC for direct `GKGSUpdateSchema` calls if immediate consistency is needed.
  - ■ **Core (Feedback & Learning Integration):** Core's `Feedback & Learning Integration` sub-system (Kafka consumer) subscribes to these topics.
  - ■ **Core's Adaptation:**
    - ■ For `PolicyUpdateSchema`: Core's `HTN Planning Engine` immediately loads the new policies, which will influence future `ActionApprovalRequests` and planning decisions.
    - ■ For `FeedbackSchema`: This data is fed into Core's internal reinforcement learning loop (if applicable) or rule refinement engine, improving its planning heuristics (e.g., learning to prefer certain tools for certain tasks, or to avoid known problematic action sequences).
    - ■ For `GKGSUpdateSchema`: Core directly applies the knowledge updates to its `Global Knowledge Graph Store`.

---

### 5. Script Kitty.Core ↔ Script Kitty.Foundry: Orchestration & Infrastructure Management

Foundry is the underlying platform that provides all the compute, MLOps, and security services necessary for Core to operate and manage other agents. Core's interaction with Foundry is typically through Foundry's Kubernetes API and MLOps tools.

- **5.1. Core to Foundry: Agent/Resource Provisioning & Orchestration**

  - **Purpose:** Core dynamically requests Foundry to provision, scale, or terminate specialized agent instances (Skills, Armory) or specific compute resources (GPUs, VMs) based on its plan.
  - **Flow & Data Schema:**
    - **Core (Agent Orchestration & Lifecycle Management):** When Core's plan requires invoking a specific Skill agent (e.g., a rarely used, resource-heavy LLM for code generation) or a custom Armory tool requiring a dedicated environment, Core needs to ensure the underlying infrastructure is ready.
    - **Core (Internal Communication Hub/Kubernetes Client):** Core doesn't directly deploy Docker containers. Instead, it interacts with Foundry's exposed APIs:
      - **Kubernetes Python Client:** Core uses the Kubernetes API client to create/update/delete custom resources (CRDs) defined by Foundry, such as `SkillModule` or `ArmoryToolRunner` CRDs.
        - Example: Core creates a `SkillModule` CR, specifying `skill_id="CodeLlama-70B"`, `replica_count=1`, `gpu_type="A100"`.
      - **Argo Workflows/Kubeflow Pipelines API:** For complex multi-step agent workflows that span multiple Skills or Armory interactions, Core might trigger a pre-defined workflow template managed by Foundry's orchestration layer (e.g., "run a complex data analysis pipeline").
  - **Communication Protocol:** Kubernetes API (RESTful over HTTPS), gRPC/REST for MLOps platform APIs (e.g., Kubeflow Pipelines API).
  - **Foundry (Kubernetes-Native Compute & Orchestration):** Foundry's custom controllers for `SkillModule` CRDs (or similar) watch for these requests. Upon detection, Foundry translates them into standard Kubernetes Deployments, Services, and Pods, ensuring the correct Docker images are pulled from the `Container Registry`, resources are allocated, and network policies are applied. Foundry dynamically scales these components using HPA/VPA.
  - **Foundry (Resource Provisioning via IaC Templates):** If a request involves external infrastructure (e.g., a specific cloud VM), Foundry's IaC orchestrator (Terraform/Ansible) handles the provisioning.
- **5.2. Foundry to Core: Infrastructure Status & Operational Metrics**

  - **Purpose:** Foundry continuously provides Core with critical information about the health, availability, and performance of the underlying infrastructure and managed agents.
  - **Flow & Data Schema:**

- **Foundry (Observability Stack):** Foundry's Prometheus, Grafana, Loki, and OpenTelemetry services collect massive amounts of data from all components, including Core, Nexus, Armory, and Skills.
- **Foundry (Internal Metrics/Log Streaming):** Foundry can expose aggregated metrics or specific log streams to Core via:
    - **Prometheus Query Language (PromQL):** Core's `Agent Orchestration` might periodically query Foundry's Prometheus for high-level metrics (e.g., `average_skills_agent_latency`, `available_gpu_resources`, `pod_failure_rate_by_agent_type`).
    - **Kafka Topics:** Foundry can publish aggregated `InfrastructureStatusSchema` (Protobuf) messages to dedicated Kafka topics (e.g., `sk_infrastructure_status`, `sk_resource_availability`). These messages might include:
        - `resource_type`
        - `available_count`
        - `total_count`
        - `health_status` (OK, DEGRADED, CRITICAL)
        - `agent_type_health` (e.g., map of Skill type to health status)
  - **Communication Protocol:** Primarily through Prometheus queries (HTTP/S) and Kafka.
  - **Core (Agent Orchestration & Lifecycle Management):** Core's `Orchestrator` consumes these infrastructure updates.
  - **Core's Adaptation:** Core uses this information to make more informed planning decisions:
    - If GPU resources are low, Core might prioritize CPU-bound tasks or defer GPU-intensive ones.
    - If a specific `SkillModule` is reported as `DEGRADED`, Core might avoid delegating tasks to it or try a different version.
    - This closed-loop feedback allows Core to adapt its operational plan to the realities of the underlying infrastructure, directly impacting its efficiency and reliability.

---

**Recap of Inter-Component Flow & Core Principles:**

- **Decoupling:** Kafka provides critical decoupling, allowing agents to operate at their own pace and preventing cascading failures. A slow Skills agent won't block Core if Core is consuming from a Kafka topic.

- **High-Fidelity Communication:** gRPC and Protobuf ensure precise, high-performance, and type-safe interaction for critical requests and responses.
- **Observability-Driven Adaptability:** Foundry's comprehensive observability stack provides the raw data that Guardian uses for learning and that Core uses for real-time operational adjustments. Every interaction leaves a digital trace.
- **Security by Design:** mTLS, network policies, secrets management, and runtime security measures are baked into every communication path, ensuring that Script Kitty's internal operations are robust against compromise.
- **Global Knowledge Graph (GKGS) as Shared Context:** While not a direct communication channel, the GKGS acts as a universally accessible shared memory. Core constantly updates and queries it, and other agents can contribute or retrieve information through Core or dedicated APIs (with Guardian oversight). This ensures all agents operate with a consistent, ever-evolving understanding of the world and Script Kitty's capabilities. For instance, Armory's digested research findings are stored in GKGS by Core, making them available to Skills for task execution, or Guardian for policy evaluation.
- **Continuous Evolution:** Every interaction generates data (logs, metrics, task results, feedback) that flows back to Foundry's data lake and Guardian's evaluation engine, fueling the continuous retraining and refinement of all Script Kitty's models and planning heuristics. This closed-loop learning across inter-component interactions is foundational to Script Kitty's AGI capabilities.

# Script Kitty.Armory

## Part 3: Script Kitty.Armory - Resource Acquisition & Tooling AI

**Core Purpose:** Script Kitty.Armory serves as the system's external interface to the vast ocean of global information and computational resources. Its fundamental mission is to dynamically acquire, meticulously process, and intelligently digest external data, ranging from raw web content to specialized research papers. Concurrently, it maintains a living, breathing registry of all available external tools and programmatic interfaces, from command-line utilities to sophisticated quantum computing APIs, ensuring Script Kitty has a dynamic and expansive operational capability. It is the architect of Script Kitty's capacity to *reach out* and *integrate* the external world into its internal cognitive and operational framework.

**Technical Stack Overview (Armory):**

- **Languages:** Predominantly Python (for web scraping, NLP, data processing, and cloud SDKs) with Go (for high-concurrency network operations and API gateways).
- **Web Automation:** Playwright (preferred for modern web apps), Selenium.
- **HTML Parsing:** BeautifulSoup4, lxml.
- **Text Extraction:** Trafilatura, boilerpy3.
- **NLU/LLM for Data Processing:** Fine-tuned LLM (e.g., Llama 3-8B), spaCy, Hugging Face Transformers.
- **Databases:** PostgreSQL (for structured tool manifests), Elasticsearch/OpenSearch (for semantic tool discovery).
- **Code Generation:** Code-focused LLM (e.g., CodeLlama, StarCoder2).
- **IaC:** Terraform, Ansible.
- **Cloud SDKs:** boto3 (AWS), google-cloud-python (GCP).
- **Quantum SDKs:** Qiskit, Cirq, PennyLane.
- **Orchestration:** Kubernetes API client, Argo Workflows.
- **Schema Definition:** Protocol Buffers (Protobuf), Pydantic.

---

**Minute Aspects & Technical Dissection:**

**1. Intelligent Web Scraping & Data Digestion Sub-system**

This sub-system is the digital equivalent of a seasoned researcher, capable of autonomously navigating the web, extracting precisely what's needed, and transforming it into actionable intelligence.

- **1.1. Research Query Interpretation & Prioritization (Armory LLM):**

  - **Purpose:** Takes a high-level `ResearchQuerySchema` from Script Kitty.Core (e.g., "Find the latest research on explainable AI in healthcare," "Gather technical specifications for the XYZ API"). It then refines this into concrete search terms,

identifies relevant domains, and potentially breaks down the query into sub-queries.

- ○ **Technical Dissection:**
  - ■ **LLM Model:** A dedicated instance of a moderately-sized LLM (e.g., Llama 3-8B), possibly fine-tuned on search query refinement and information retrieval tasks. This LLM is served by `vLLM` or `Ollama` for efficient inference.
  - ■ **Prompt Engineering:** The LLM receives a system prompt instructing it to act as a "research strategist." It's fed the `ResearchQuerySchema` (including keywords, desired output format, constraints like publication date, source type, etc.) and its output is a structured `WebSearchPlanSchema` (Protobuf), detailing search engines to use, initial keywords, desired document types (PDF, HTML), and a list of potential domains to prioritize.
  - ■ **Search Engine Integration:** Uses `requests` or dedicated Python libraries to interact with major search engine APIs (e.g., Google Custom Search API, Bing Web Search API, specialized academic search APIs like Semantic Scholar or PubMed). It parses JSON responses to extract initial URLs.
  - ■ **Query Expansion & Refinement:** Based on initial search results and feedback from the scraper, the LLM can iteratively refine search queries or generate new ones to improve relevance and coverage.
  - ■ **Result Prioritization:** Implements ranking algorithms (e.g., TF-IDF, BM25, or even semantic similarity using embeddings) to prioritize search results based on estimated relevance to the original `ResearchQuerySchema`, considering factors like domain authority, recency, and keyword density.
- **1.2. Dynamic Web Navigation & Content Extraction (Scraper Agents):**

  - ○ **Purpose:** Executes the `WebSearchPlanSchema` by navigating web pages, bypassing common anti-scraping measures, and extracting raw content.
  - ○ **Technical Dissection:**
    - ■ **Headless Browser Automation (Playwright/Selenium):**
      - ■ `Playwright` (preferred for its modern API, speed, and ability to handle JavaScript-heavy SPAs) or `Selenium` (for broader browser compatibility and legacy support) is used.
      - ■ `Chromium` or `Firefox` instances run in headless mode within isolated Docker containers.
      - ■ **Proxy Rotation:** Integrates with a proxy pool service (e.g., using `luminati-proxy` or custom solution with a large pool of residential/datacenter proxies) to prevent IP blocking. This is managed by a `ProxyManager` microservice.

- **User-Agent Rotation:** Employs a dynamic rotation of diverse user-agent strings to mimic legitimate browser traffic.
- **CAPTCHA Bypass Integration (Optional/Third-Party):** For sites with CAPTCHA, integrates with third-party CAPTCHA solving services (e.g., 2Captcha, Anti-Captcha) via their APIs, with fallbacks and cost monitoring.
- **DOM Interaction:** Uses browser automation APIs to click buttons, fill forms, scroll to load dynamic content, and navigate paginated results, simulating human Browse behavior.
- **JavaScript Execution:** Ensures that all necessary JavaScript on the page is executed to render the full content before extraction.
- **Screenshot/PDF Generation:** Can generate screenshots or PDFs of pages for auditability or specific content types.
- **HTML Parsing & Text Extraction (BeautifulSoup4/lxml, Trafilatura/boilerpy3):**
  - `BeautifulSoup4` (Python) or `lxml` (for speed with large XML/HTML documents) is used to parse the HTML DOM.
  - `CSS selectors` or `XPath` are used for robust and precise extraction of specific elements (e.g., article content, headings, lists, tables).
  - `Trafilatura` or `boilerpy3` is employed to intelligently identify and extract the main content/article from a web page, stripping boilerplate (headers, footers, ads, navigation), ensuring clean text for NLP.
  - **Metadata Extraction:** Extracts key metadata like title, author, publication date, and canonical URL.
- **Error Handling & Retries:** Implements sophisticated error handling for network errors, timeouts, invalid selectors, and anti-bot measures. Includes configurable retry strategies with exponential backoff and dynamic delays.
- **Rate Limiting:** Enforces per-domain and global rate limits to avoid overwhelming target servers and to maintain good web citizenship.

- **1.3. Information Processing & Structuring (NLP Pipeline):**

  - **Purpose:** Transforms raw, extracted text into structured, usable data, including summaries, entities, and relationships.
  - **Technical Dissection:**
    - **Pre-processing & Cleaning:** Further cleans the extracted text: removes extraneous characters, normalizes encoding, corrects minor OCR errors if dealing with PDFs, and performs sentence segmentation.
    - **Named Entity Recognition (NER) (spaCy/Hugging Face Transformers):**

- Utilizes pre-trained or fine-tuned NER models (e.g., `spaCy`'s `en_core_web_lg` or transformer-based models from `Hugging Face` like `bert-base-NER`) to identify and classify entities (e.g., persons, organizations, locations, dates, technical terms, specific API names, command flags).
- **Custom Entity Models:** For domain-specific information (e.g., software vulnerabilities, specific scientific jargon), custom NER models are trained using annotated data, potentially with active learning.
- **Relation Extraction (RE):** Identifies semantic relationships between extracted entities (e.g., "Python (language) *uses* Pandas (library)," "Kubernetes (platform) *manages* Pods (resource)"). Can use rule-based patterns or transformer-based RE models.
- **Summarization (Hugging Face Transformers):** Employs abstractive or extractive summarization models (e.g., `bart-large-cnn`, `t5-small`) to condense lengthy articles into concise summaries, tailored to the `ResearchQuerySchema`'s specified length and focus.
- **Keyword/Topic Extraction:** Uses techniques like TextRank, LDA, or BERT-based topic modeling to identify key themes and keywords from the processed text.
- **Sentiment Analysis (Optional):** For specific research queries involving public opinion or reviews, a sentiment analysis model can be applied.
- **Data Validation & Schema Enforcement (Pydantic):** The processed data is validated against a `ProcessedDataSchema` (Protobuf/Pydantic model) to ensure consistency and adherence to predefined output structures (e.g., list of `Entity` objects, `Summary` string, `Keywords` array). This ensures that Core receives data in a predictable format.
- **Semantic Embedding Generation:** Generates vector embeddings for extracted text segments, summaries, and entities using models like Sentence-BERT (`sbert`) or OpenAI's embedding models (via API). These embeddings are then stored in the Global Knowledge Graph for semantic search and RAG by other agents.

**2. Dynamic Tool Registry & Discovery Sub-system**

This sub-system acts as Script Kitty's comprehensive inventory of all external capabilities it can leverage, providing a discoverable and structured interface for the Core.

- **2.1. Tool Manifest Database:**

  - **Purpose:** Stores and manages structured metadata for every tool accessible to Script Kitty, making it discoverable and understandable by Core's planning engine.
  - **Technical Dissection:**

- **Database:** PostgreSQL is chosen for its robust relational capabilities, transactional integrity, and strong support for JSONB data types (for flexible `capabilities` and `parameters` schemas).
- **`ToolManifestSchema` (Protobuf):** This is the heart of the registry. It's a rigorously defined schema encompassing:
    - `tool_id` (UUID), `tool_name`, `description` (detailed natural language description of what the tool does).
    - `capabilities` (an array of semantic tags/keywords, e.g., "web_scraping", "image_processing", "cloud_provisioning", "code_execution", "quantum_simulation"). These are critical for Core's semantic search.
    - `input_parameters` (a JSON Schema or Protobuf message detailing expected arguments, their types, validation rules, and descriptions).
    - `output_schema` (a JSON Schema or Protobuf message describing the structure of the tool's return value).
    - `required_environment` (e.g., `python_3.10`, `docker_image:some_tool_v1.0`, `aws_cli_installed`).
    - `adapter_type` (e.g., "CLI_wrapper", "Python_library", "REST_API_client", "gRPC_client").
    - `security_context` (required permissions, isolation level).
    - `usage_cost_model` (e.g., API calls, compute time, flat fee – for resource budgeting).
- **Indexing:** Critical fields like `capabilities`, `tool_name`, and `description` are indexed for fast retrieval.
- **2.2. Semantic Tool Search & Discovery Engine (Elasticsearch/OpenSearch):**

    - **Purpose:** Allows Script Kitty.Core to query the registry using natural language descriptions of desired functionalities and retrieve the most relevant tools.
    - **Technical Dissection:**
        - **Search Engine:** Elasticsearch or OpenSearch (open-source fork) is used for its powerful full-text search, semantic search capabilities, and scalability.
        - **Index Structure:** Tool manifests are indexed with fields including `description`, `capabilities`, and `tool_name`.
        - **Embedding Search:** Vector embeddings of the `description` and `capabilities` fields are generated (e.g., using `sentence-transformers` or `OpenAI` embeddings) and stored in Elasticsearch using its `dense_vector` type. This enables k-nearest neighbor (KNN) or Approximate Nearest Neighbor (ANN) search for semantic similarity.

- **Hybrid Search:** Combines keyword search (e.g., BM25) with semantic vector search for highly relevant results.
- **Query Translation:** Script Kitty.Core sends a natural language query describing desired functionality to Armory. Armory's internal LLM or a small fine-tuned model translates this into an Elasticsearch query (combining keywords and vector searches).
- **Filtering & Ranking:** Search results are filtered by `required_environment` and `security_context` and ranked based on relevance, cost, and availability.

**3. Templated API/CLI Wrapper Generation (Human-Verified) Sub-system**

This highly innovative sub-system automates the tedious and error-prone process of integrating new external tools by generating their programmatic wrappers. This is a critical enabler for Script Kitty's extensibility.

- **3.1. Wrapper Code Generation (Code LLM):**

  - **Purpose:** Takes a `ToolManifestSchema` (for a new or updated tool) and automatically generates boilerplate code for a Python or Go wrapper, allowing Script Kitty.Skills to invoke the tool.
  - **Technical Dissection:**
    - **Code-Focused LLM:** Utilizes a specialized LLM (e.g., CodeLlama-7B/34B, StarCoder2, DeepSeek Coder) specifically fine-tuned on a vast corpus of API client libraries, CLI wrappers, and common programming patterns. This LLM is served by `vLLM` or `Text Generation Inference`.
    - **Prompt Engineering:** The LLM receives a detailed prompt including the `ToolManifestSchema`, the desired `adapter_type` (e.g., Python `requests` wrapper, `subprocess` CLI call), and example invocation patterns. The prompt instructs the LLM to generate code that:
      - Parses input parameters according to `input_parameters`.
      - Constructs HTTP requests (for REST APIs) or CLI commands correctly.
      - Handles authentication (e.g., API keys from Vault).
      - Parses tool output into the `output_schema`.
      - Includes basic error handling and logging.
      - Adheres to predefined coding style guidelines.
    - **Templating Engine (Jinja2):** For structural consistency and security, the LLM-generated code is often injected into pre-defined, secure wrapper templates. For example, a template might define the class structure, exception handling, and input validation, while the LLM fills in the specific API call logic.

- ■ **Security Context Integration:** The generated wrapper code is designed to respect the `security_context` (e.g., ensuring secrets are accessed via Vault, not hardcoded; enforcing least privilege access).
- ● **3.2. Static Analysis & Security/Quality Gate:**

  - ○ **Purpose:** Automatically scans the LLM-generated code for syntax errors, potential security vulnerabilities, and quality issues *before* human review.
  - ○ **Technical Dissection:**
    - ■ **Syntax & Linting:** Uses language-specific linters (e.g., `flake8`, `pylint` for Python; `golint`, `gofmt` for Go) to ensure syntactical correctness and adherence to style guides.
    - ■ **Static Application Security Testing (SAST):**
      - ■ `Bandit` (for Python) or `semgrep` (multi-language) are used to detect common security vulnerabilities (e.g., SQL injection, insecure deserialization, weak cryptography, hardcoded credentials).
      - ■ `SonarQube Community Edition` provides broader static analysis for code quality, maintainability, and architectural issues.
    - ■ **Dependency Scanning:** `Trivy` or `Snyk Open Source` scan the `requirements.txt` or `go.mod` files for known vulnerabilities in third-party libraries.
    - ■ **Automated Unit Test Generation (Optional/Future):** The code LLM could also generate basic unit tests for the wrapper, and these tests are run by a testing framework (e.g., `pytest`).
    - ■ **Policy Enforcement:** The generated code is evaluated against a set of predefined policies using `Open Policy Agent (OPA)` rules (e.g., "no direct filesystem access," "only allow HTTPs connections").
    - ■ **Report Generation:** Generates a detailed report of all findings, including severity levels, which is then presented to the human reviewer.
- ● **3.3. Human Review & Approval Interface:**

  - ○ **Purpose:** Provides a critical human-in-the-loop step, allowing human operators to review, modify, and explicitly approve LLM-generated tool wrappers before they are integrated into the system. This ensures safety and correctness.
  - ○ **Technical Dissection:**
    - ■ **Web Interface (React/Next.js Frontend, FastAPI Backend):** A dedicated web application serves as the review portal.
    - ■ **Diff View:** Presents a clear, side-by-side diff between the generated code and any previous version, highlighting changes.
    - ■ **Automated Analysis Report Display:** Integrates the output from the static analysis tools directly into the UI, making it easy for reviewers to see identified issues.

- **Annotation & Editing:** Allows human reviewers to directly edit the generated code within the interface, or add annotations and comments for the LLM's improvement feedback loop.
- **Approval Workflow:** Implements a clear approval button. Once approved, the code is committed to a secure Git repository, triggering a CI/CD pipeline for integration and deployment via Foundry.
- **Reject/Feedback Loop:** A "Reject" option allows reviewers to provide detailed feedback, which is then fed back into the LLM's fine-tuning data or prompts for future generation attempts, continuously improving its accuracy and safety.
- **Audit Trail:** All review actions, modifications, and approvals are logged for auditability and compliance.

## 4. Resource Provisioning via IaC Templates Sub-system

This sub-system enables Script Kitty to dynamically provision the necessary computational and data resources it needs to execute tasks, acting as its self-provisioning infrastructure manager.

- **4.1. IaC Template Library:**

  - **Purpose:** A curated repository of pre-approved, secure, and parameterized Infrastructure as Code (IaC) templates for various cloud resources.
  - **Technical Dissection:**
    - **Template Storage:** IaC blueprints (e.g., Terraform HCL, Ansible Playbooks, Kubernetes YAMLs) are stored in a version-controlled Git repository (e.g., GitLab, GitHub Enterprise) that is accessible only to Armory.
    - **Parameterization:** Templates are heavily parameterized to allow dynamic configuration (e.g., instance type, region, storage size, security groups, quantum resource allocation) based on Core's request.
    - **Security & Compliance:** All templates undergo rigorous security review by Guardian and human experts, ensuring they adhere to organizational security policies and best practices (e.g., least privilege, encryption at rest/in transit).
    - **Versioning:** Templates are versioned, allowing for rollbacks to previous stable configurations.
- **4.2. IaC Orchestrator & Cloud/Quantum SDK Integrator:**

  - **Purpose:** Interacts with cloud provider APIs and quantum computing services to provision, modify, or de-provision resources based on Core's `ResourceProvisioningRequestSchema`.
  - **Technical Dissection:**
    - **Terraform CLI Automation:** Armory invokes the `terraform` CLI within a secure, sandboxed environment. It dynamically generates Terraform

variable files (`.tfvars`) based on the `ResourceProvisioningRequestSchema` received from Core. It executes `terraform init`, `terraform plan`, and `terraform apply`.

- **Ansible Automation:** For configuration management on provisioned VMs (e.g., installing specific software, configuring network interfaces), Armory generates and executes Ansible Playbooks.
- **Kubernetes API Client:** For provisioning ephemeral Kubernetes resources (e.g., a GPU-enabled Pod for a specific ML task, a custom PVC for storage), Armory directly interacts with the Kubernetes API using the `kubernetes` Python client library. This allows for fine-grained, dynamic resource allocation within the existing cluster.
- **Cloud Provider SDKs (boto3, google-cloud-python):** For scenarios where full IaC frameworks are overkill or for direct programmatic interaction with specific services (e.g., managing S3 buckets, triggering serverless functions, interacting with specific ML services), Armory uses official cloud SDKs. Authentication uses temporary credentials issued by Vault.
- **Quantum Computing SDKs (Qiskit, Cirq, PennyLane):** Armory integrates with these SDKs to programmatically access quantum computing hardware or simulators (e.g., IBM Quantum Experience, Google Quantum AI, Xanadu's PennyLane cloud). It manages access keys (from Vault), job submission, monitoring, and result retrieval. It converts Core's abstract quantum task requests into platform-specific quantum circuits or programs.
- **Idempotency & State Management:** Terraform inherently provides idempotency, ensuring consistent deployments. Armory carefully manages Terraform state files, often storing them in a secure backend (e.g., S3/GCS bucket with state locking).
- **Monitoring & Rollback:** Monitors the provisioning process for errors and timeouts. If provisioning fails, it triggers automated rollback mechanisms (e.g., `terraform destroy`) and reports failures to Guardian.
- **Cost Monitoring & Optimization:** Integrates with cloud cost management APIs to track resource consumption and provides feedback to Core and Guardian for cost optimization during planning.

---

**Operational Excellence & Extensibility (Armory-Specific):**

- **Observability (Integrated with Foundry):**

  - **Metrics:** Detailed metrics on `scraping_success_rate`, `tool_invocation_latency`, `wrapper_generation_time`,

`resource_provisioning_duration`, `cloud_api_call_rate`, and `quantum_job_status` are collected via Prometheus.
- **Logging:** Comprehensive structured logs for all scraping activities, tool registry updates, and IaC operations, aggregated by Loki/Elasticsearch. Includes granular details like URLs visited, parsing errors, generated code, and IaC command outputs.
- **Tracing:** OpenTelemetry traces capture the full lifecycle of a research query or a tool integration request, spanning multiple Armory sub-components and external API calls.
- **Alerting:** Alerts triggered for `Scraping_Blocked_High`, `Tool_Wrapper_Generation_Failure`, `Resource_Provisioning_Timeout`, `Quantum_Job_Failed`.

- **Scalability:**

  - **Scraper Agents:** The `Scraper Agents` (running headless browsers) are deployed as Kubernetes Jobs or Deployments, horizontally scaling based on the volume of research queries. Each agent runs within a tightly resource-limited and network-isolated container.
  - **IaC Orchestrator:** Designed to handle concurrent resource provisioning requests by leveraging job queues (e.g., Kafka, Redis queues) and managing multiple concurrent Terraform/Ansible executions.
  - **LLM Inference:** The dedicated LLM inference services for NLU and code generation are horizontally scalable, utilizing GPU/CPU clusters.

- **Continuous Improvement & Extensibility (Via Script Kitty.Foundry & Guardian):**

  - **Feedback Loops:** Human feedback from the Wrapper Review Interface and performance metrics from the IaC Orchestrator (e.g., provisioning success/failure, duration) are fed back to Guardian and Foundry. This data is used to fine-tune the Code LLM for better wrapper generation and to refine IaC templates.
  - **Automated Tool Onboarding (Future):** Armory aims to automate the initial ingestion of tool documentation (e.g., OpenAPI specs, CLI `--help` outputs) to automatically generate initial `ToolManifestSchemas` (requiring human verification).
  - **New IaC Templates:** The library of IaC templates is continuously expanded to support new cloud services or quantum computing backends as they become available.
  - **Dynamic LLM Fine-tuning:** The LLMs within Armory are continuously fine-tuned with new, annotated data generated from successful research queries and human-corrected wrapper code, enhancing its ability to understand information and generate correct code.

# Armory Interactions

# Part 7: Inter-Component Interactions - The Orchestrated Symphony

This section meticulously dissects the communication channels, data flows, and collaborative protocols that bind Script Kitty's specialized components into a cohesive, intelligent whole.

---

## Section 7.1: Script Kitty.Armory - The Nexus of Resources and Tools

Script Kitty.Armory, the system's specialized AI for acquiring, processing, and digesting external information, managing dynamic tool registries, and provisioning external resources, is deeply intertwined with every other facet of Script Kitty. Its interactions are critical for equipping the system with the capabilities and knowledge it needs to solve complex problems.

### 7.1.1. Armory ↔ Script Kitty.Core: The Strategic Supply Chain

This is the primary command-and-control channel for Armory. Core dictates the research and tool acquisition needs, and Armory fulfills these requests, providing the raw materials and operational capabilities for Core's plans.

- **7.1.1.1. Core to Armory: Research & Tool Acquisition Requests**

  - **Purpose:** Script Kitty.Core, during its `Goal Decomposition & HTN Planning` phase, identifies a need for external information (e.g., "research X topic," "find an API for Y functionality") or requires a specific external tool/resource. Core then issues a highly structured request to Armory.
  - **Communication Protocol:** gRPC (Remote Procedure Call) for its high performance, bidirectional streaming capabilities, and strict schema enforcement.
  - **Data Schema:** A Protobuf-defined `ArmoryRequestSchema` is used.
    - `request_id`: Unique identifier for the request, enabling traceability and correlation with Core's ongoing plan.
    - `task_id`: Inherited from the overarching `ScriptKittyTaskSchema` managed by Core, linking the Armory request to a specific planning step.
    - `request_type`: An enum (e.g., `RESEARCH_QUERY`, `TOOL_ACQUISITION`, `IAC_PROVISIONING`, `WRAPPER_GENERATION`).
    - `query`: Natural language string or structured JSON (e.g., for research, "summarize the latest advancements in quantum computing security," or for tool, "find a Python library for image segmentation on medical scans").
    - `parameters`: A JSON object for specific instructions (e.g., `depth=3` for web scraping, `api_endpoint_pattern="google.*API"` for tool acquisition).

- **required_output_schema**: (Optional) A Protobuf or JSON schema indicating the desired format for the output (e.g., "return a list of academic papers with title, author, abstract, and URL"). This guides Armory's data digestion.
    - **deadline**: A timestamp indicating when the result is needed by Core for its planning.
- **Microscopic Detail:**
    - **Request Validation:** Armory's gRPC server rigorously validates incoming `ArmoryRequestSchema` against its definition. Invalid requests are immediately rejected with specific error codes (e.g., `INVALID_ARGUMENT`).
    - **Queueing:** Upon reception, requests are often placed into an internal, prioritized message queue (e.g., backed by Kafka or RabbitMQ) within Armory to handle bursts of requests from Core and ensure ordered processing. Critical requests might bypass the queue.
    - **Contextual Inheritance:** Armory internal processes retain `request_id` and `task_id` throughout their execution, ensuring all intermediate data and final results are correctly attributed back to Core's original plan.

- **7.1.1.2. Armory to Core: Research & Tool Acquisition Results / Status Updates**

    - **Purpose:** Armory reports the completion of requested tasks, provides acquired data, tool manifests, or resource provisioning confirmations. It also sends periodic status updates for long-running operations.
    - **Communication Protocol:** gRPC stream (for continuous updates) or gRPC unary call (for final results).
    - **Data Schema:** Protobuf-defined `ArmoryResponseSchema` and `ArmoryStatusUpdateSchema`.
        - **ArmoryResponseSchema:**
            - `request_id`: Matches the original `ArmoryRequestSchema`.
            - `task_id`: Matches the original `ScriptKittyTaskSchema`.
            - `status`: Enum (e.g., `SUCCESS`, `FAILURE`, `PARTIAL_SUCCESS`).
            - `result_type`: Enum (e.g., `RESEARCH_DATA`, `TOOL_MANIFEST`, `IAC_CONFIRMATION`, `WRAPPER_CODE`).
            - `data`: A byte array or string containing serialized result data (e.g., JSON for research, Protobuf for tool manifest, code string for wrapper). This data is often a reference to a larger object stored in Foundry's object storage (e.g., S3 URL) to avoid sending large payloads directly via gRPC.
            - `error_details`: If `status` is `FAILURE`, detailed error messages, stack traces, and relevant logs.

- - - `metadata`: Optional key-value pairs for additional context (e.g., `source_url` for research, `estimated_cost` for resource provisioning).
    - **`ArmoryStatusUpdateSchema`:**
      - `request_id`, `task_id`: As above.
      - `progress_percentage`: Integer (0-100).
      - `message`: Human-readable status (e.g., "Web scraping 50% complete," "Generating tool wrapper...").
      - `estimated_time_remaining`: Optional duration.
      - `sub_task_status`: Details on current sub-task within Armory's execution.
  - **Microscopic Detail:**
    - **Idempotency:** Core's system for receiving Armory responses is designed to be idempotent, gracefully handling duplicate `ArmoryResponseSchema` messages (e.g., due to network retries) to prevent erroneous re-processing.
    - **Callback Mechanism:** Core registers a callback or sets up a dedicated listener for responses linked to its `request_id` and `task_id`, allowing it to resume planning upon Armory's completion.
    - **Result Persistence:** Before sending, Armory ensures critical results (e.g., parsed research data, generated code) are durably stored in Foundry's object storage and DVC, with the `ArmoryResponseSchema.data` field containing only a secure, versioned reference (e.g., a signed S3 URL or DVC path).
- **7.1.1.3. Core to Armory: Tool Wrapper Generation Requests (Specialized)**

  - **Purpose:** Core might require a custom wrapper for a newly acquired or abstract tool, or a specialized CLI command. It sends detailed specifications to Armory for code generation.
  - **Communication Protocol:** gRPC (Unary or Stream).
  - **Data Schema:** An `ArmoryWrapperGenerationRequestSchema` extends `ArmoryRequestSchema`.
    - `tool_description`: Natural language description of the tool's purpose and functionality.
    - `api_spec`: (Optional) OpenAPI/Swagger spec, or a simplified JSON schema for the tool's API/CLI.
    - `target_language`: (e.g., "python", "javascript", "bash").
    - `input_schema`: Protobuf/JSON schema for the wrapper's expected inputs.
    - `output_schema`: Protobuf/JSON schema for the wrapper's expected outputs.

- - ■ `security_constraints`: (e.g., `network_isolated=true`, `no_internet_access`).
    - ■ `template_id`: (Optional) Reference to a specific pre-approved code generation template.
  - ○ **Microscopic Detail:**
    - ■ **LLM Prompt Injection:** Armory's internal `Templated API/CLI Wrapper Generation` sub-component uses the `tool_description`, `api_spec`, and `input/output_schema` to construct highly specific prompts for its code-focused LLM (e.g., CodeLlama).
    - ■ **Human-in-the-Loop Integration:** For critical or novel wrapper generations, Armory's response includes a flag (`requires_human_review=true`) and details on how Guardian's policy enforcement will handle the review process before final deployment.
- ● **7.1.1.4. Armory to Core: Tool Manifest Provisioning**

  - ○ **Purpose:** After acquiring or generating a tool, Armory registers its capabilities and usage instructions, making it discoverable for Core's planning and Skills' execution.
  - ○ **Communication Protocol:** gRPC.
  - ○ **Data Schema:** A `ToolManifestSchema` (Protobuf) is sent to Core.
    - ■ `tool_id`: Unique identifier (e.g., `web_scraper_v2`, `image_segmentation_api`).
    - ■ `name`, `description`: Human-readable details.
    - ■ `capabilities`: A list of semantic tags or natural language descriptions (e.g., "image processing," "text summarization," "cloud resource provisioning"). Crucial for Core's semantic search.
    - ■ `input_schema`, `output_schema`: Detailed Protobuf/JSON schemas for the tool's inputs and outputs.
    - ■ `execution_environment`: (e.g., `PYTHON_SANDBOX`, `CLI_CONTAINER`, `EXTERNAL_API_CALL`).
    - ■ `adapter_details`: Specifics on how to invoke the tool (e.g., API endpoint, CLI command arguments, Python module path).
    - ■ `resource_requirements`: (e.g., `cpu_cores=2`, `memory_gb=4`, `gpu_type=T4`).
    - ■ `security_profile`: (e.g., `internet_access=true`, `disk_write_allowed=false`).
    - ■ `version`: Semantic version of the tool/wrapper.
    - ■ `status`: (e.g., `ACTIVE`, `DEPRECATED`, `UNDER_REVIEW`).
  - ○ **Microscopic Detail:**

- **Knowledge Graph Update:** Core (specifically its `Global Knowledge Graph (GKGS)` component) ingests these `ToolManifestSchema` objects. The `capabilities` field is critical here, as Core's semantic search over the GKGS uses embeddings of these descriptions to match tasks to available tools.
- **Real-time Availability:** Once the `ToolManifestSchema` is acknowledged by Core's GKGS, the tool becomes immediately discoverable and usable by Core's planning engine and potentially by Skills agents for direct execution.
- **Version Management:** Core stores multiple versions of `ToolManifestSchema` for the same `tool_id`, allowing Core's planner to select specific versions or fall back to older, stable ones if newer versions encounter issues.

**7.1.2. Armory ↔ Script Kitty.Nexus: Indirect User Interface for Resource Acquisition**

Direct interaction is minimal, as Nexus is the primary user interface. Armory's interactions with Nexus are primarily indirect, flowing through Core, which then translates updates for the user.

- **7.1.2.1. Armory to Core (then to Nexus): Status Updates for User Visibility**

  - **Purpose:** When Armory is performing a long-running task (e.g., extensive web scraping, complex resource provisioning), it sends `ArmoryStatusUpdateSchema` messages to Core. Core then translates these into user-friendly `ScriptKittyStatusUpdateSchema` messages for Nexus to display to the user.
  - **Communication Flow:** Armory → gRPC → Core → gRPC → Nexus.
  - **Microscopic Detail:**
    - **Core's Translation Layer:** Core contains a specialized `TaskStatusTranslator` that takes `ArmoryStatusUpdateSchema` and converts it into a more abstract, user-facing `ScriptKittyStatusUpdateSchema`. This abstraction ensures that users don't see raw technical details from Armory but rather high-level progress updates (e.g., "Armory is researching X topic," instead of "Armory running Playwright browser automation on URL Y").
    - **Nexus's Persona Layer:** Nexus's `Response Synthesis & Persona Layer` receives these updates from Core and synthesizes them into conversational text, potentially adding persona-specific flair (e.g., "My esteemed colleague, Armory, is diligently scouring the digital expanse for that information you requested...").
- **7.1.2.2. Armory to Core (then to Nexus): Final Results for User Consumption**

- ○ **Purpose:** The final outcome of Armory's work (e.g., summarized research, confirmation of resource provisioning) is relayed to the user through Core and Nexus.
- ○ **Communication Flow:** Armory → gRPC `ArmoryResponseSchema` → Core → gRPC `ScriptKittyTaskResultSchema` → Nexus.
- ○ **Microscopic Detail:**
  - ■ **Core's Result Aggregation:** If Armory's task is part of a larger plan, Core might wait for results from multiple agents (Armory, Skills) before aggregating them into a single `ScriptKittyTaskResultSchema` that summarizes the overall task completion.
  - ■ **Content Sanitization (Core/Nexus):** Before presenting research results to the user, Core or Nexus (via Guardian's influence) might perform a final content sanitization pass to remove any detected PII, harmful content, or irrelevant noise that Armory might have inadvertently collected.
  - ■ **Interactive Presentation:** If Armory's result is a structured dataset, Nexus might instruct the frontend UI (if applicable) to render it as an interactive table or graph, leveraging its `suggested_actions` field.

### 7.1.3. Armory ↔ Script Kitty.Skills: The Tool Provider and Consumer of Execution Feedback

Armory provides the tools and resources for Skills to execute, while Skills, in turn, provides valuable feedback on the real-world performance of these tools.

- ● **7.1.3.1. Armory to Skills (via Core's Orchestration): Tool Provisioning for Execution**

  - ○ **Purpose:** While Core orchestrates, Armory is responsible for making the *actual* executable tool (e.g., a Python wrapper, a CLI binary) available to the relevant Skills agent instance.
  - ○ **Communication Protocol:** Primarily via Foundry (Kubernetes Persistent Volumes/Mounts, Container Registry). Core orchestrates the deployment.
  - ○ **Data Flow:**
    - ■ Armory generates or acquires a tool/wrapper (`WRAPPER_CODE` or `TOOL_BINARY`).
    - ■ Armory stores this artifact securely in Foundry's `Container Registry` (if it's a new Docker image) or `Object Storage` (if it's a script/binary).
    - ■ Armory updates its `ToolManifestSchema` and sends it to Core.
    - ■ Core, during planning, identifies a `SkillModule` (a type of `ScriptKittyAgent`) that needs this tool.
    - ■ Core (via Foundry) orchestrates the deployment/update of the `SkillModule` pod.

- Foundry's Kubernetes-native capabilities (e.g., mounting a Persistent Volume Claim (PVC) from Object Storage, or pulling the new Docker image containing the tool) make the tool available inside the Skills agent's sandbox.

- ○ **Microscopic Detail:**
  - **Secure Mounting:** For sensitive tools or dynamically generated wrappers, Foundry can create ephemeral PVCs that securely mount the tool's execution environment into the Skills agent's sandboxed container. These PVCs are configured with strict permissions and limited lifespan.
  - **Image Updates:** If a tool requires an updated environment or is a new core capability, Armory triggers the build of a new Docker image for the relevant `SkillModule` via Foundry's CI/CD. This image is then rolled out to Skills instances.
  - **Resource Templates:** Armory's `Resource Provisioning via IaC Templates` might provide Skills with access to specialized compute resources (e.g., a specific GPU type or a quantum computing simulator instance) that are required for a particular tool's execution. This is orchestrated by Foundry based on Armory's output.

- **7.1.3.2. Skills to Armory (via Core/Guardian): Tool Execution Feedback**

  - ○ **Purpose:** Skills agents, after executing a tool provided by Armory, provide feedback on its real-world performance, success rate, and any errors encountered. This is crucial for Armory's `Feedback & Learning Integration` and for refining its `Dynamic Tool Registry`.
  - ○ **Communication Flow:** Skills → gRPC `ExecutionResultSchema` → Core → gRPC `TaskFeedbackSchema` → Guardian (for policy/evaluation) → Kafka/gRPC `ToolPerformanceMetrics` → Armory.
  - ○ **Microscopic Detail:**
    - **Granular Metrics:** Skills agents don't just report success/failure. They provide granular metrics such as `execution_duration`, `resource_consumption` (CPU, memory), `output_size`, and specific `error_codes` or `exceptions`.
    - **Correlation:** Core ensures `tool_id` and `version` are always included in the feedback, allowing Armory to correlate feedback with specific tool iterations.
    - **Guardian's Role:** Guardian acts as an intermediary, processing this feedback for `Real-time Performance Monitoring & Anomaly Detection` and for `Automated Evaluation & Training Trigger` (which might prompt Armory to re-evaluate a tool). Guardian might aggregate feedback from multiple Skills executions before forwarding it to Armory.
    - **Armory's Adaptation:** Armory uses this aggregated feedback to:

- **Rank Tools:** Update its internal ranking of tools based on observed success rates and efficiency.
- **Identify Flaky Tools:** Flag tools that frequently fail or perform inconsistently.
- **Refine Wrapper Generation:** Use error patterns from executed wrappers to refine the prompts given to its code-focused LLM, generating more robust wrappers in the future.
- **Trigger Deprecation:** Automatically mark tools as `DEPRECATED` in the `Dynamic Tool Registry` if they consistently perform poorly.

**7.1.4. Armory ↔ Script Kitty.Guardian: Oversight and Ethical Integration**

Guardian plays a vital role in ensuring Armory's actions align with safety, ethical, and operational policies, particularly concerning resource acquisition and code generation.

- **7.1.4.1. Armory to Guardian: Policy Pre-check & Approval Requests**

  - **Purpose:** For critical or sensitive operations (e.g., web scraping highly sensitive domains, provisioning expensive cloud resources, deploying newly generated tool wrappers), Armory sends a pre-check request to Guardian's `Policy Enforcement Engine`.
  - **Communication Protocol:** gRPC.
  - **Data Schema:** A `PolicyCheckRequestSchema` (Protobuf) containing:
    - `action_type`: (e.g., `WEB_SCRAPE`, `RESOURCE_PROVISION`, `CODE_DEPLOY`).
    - `payload`: Specific details about the action (e.g., `target_url`, `iac_template_id`, `generated_code_hash`).
    - `originating_agent`: "Armory".
    - `risk_score`: (Optional) Armory's own internal initial risk assessment.
    - `request_id`, `task_id`: For traceability.
  - **Microscopic Detail:**
    - **OPA (Open Policy Agent) Integration:** Guardian's `Policy Enforcement Engine` uses OPA. Armory's `PolicyCheckRequestSchema` payload is translated into a JSON input for OPA's Rego policy engine. OPA evaluates the input against defined policies (e.g., "no scraping domains in blacklist," "resource cost must be below X without human approval," "generated code must pass static analysis before deployment").
    - **Human-in-the-Loop Trigger:** If OPA policies require human review (e.g., high-cost resource provisioning), Guardian flags the request in its Human

`Oversight & Feedback Interface`, pausing Armory's operation until explicit approval/rejection is received.

- ■ **Pre-computed Rules:** For common, low-risk actions, Armory might have cached pre-approved rules from Guardian to reduce latency.

- ● **7.1.4.2. Guardian to Armory: Policy Decisions & Enforcement**

  - ○ **Purpose:** Guardian responds to Armory's policy check requests with an `ALLOW` or `DENY` decision, potentially with modifications or specific instructions.
  - ○ **Communication Protocol:** gRPC.
  - ○ **Data Schema:** A `PolicyDecisionSchema` (Protobuf) containing:
    - ■ `request_id`, `task_id`.
    - ■ `decision`: Enum (`ALLOW`, `DENY`, `REQUIRES_HUMAN_APPROVAL`).
    - ■ `reasons`: A list of strings explaining the decision (e.g., "Violates PII protection policy," "Exceeds budget threshold").
    - ■ `recommended_actions`: (e.g., "Redact PII from output," "Resubmit with lower resource request").
  - ○ **Microscopic Detail:**
    - ■ **Atomic Decision:** Armory is designed to halt execution immediately upon receiving a `DENY` decision.
    - ■ **Actionable Feedback:** If `DENY` or `REQUIRES_HUMAN_APPROVAL`, Armory uses the `reasons` and `recommended_actions` to inform Core and potentially the user about the policy violation, enabling corrective measures.
    - ■ **Callback Mechanism:** Armory maintains a callback mapping `request_id` to the ongoing operation, allowing it to resume or terminate based on Guardian's decision.

- ● **7.1.4.3. Guardian to Armory: Performance Monitoring & Quality Feedback**

  - ○ **Purpose:** Guardian's `Real-time Performance Monitoring & Anomaly Detection` (from Skills execution data) and `Automated Evaluation` pipelines provide feedback to Armory on the effectiveness and quality of the tools and data it provides.
  - ○ **Communication Protocol:** Kafka stream (for aggregated, asynchronous feedback).
  - ○ **Data Schema:** A `ToolQualityFeedbackSchema` (Protobuf) aggregated by Guardian from Skills' execution results.
    - ■ `tool_id`, `version`: Specific tool being evaluated.
    - ■ `success_rate_24h`: Rolling average of successful executions.
    - ■ `avg_latency_ms`: Average execution time.
    - ■ `error_types_distribution`: Distribution of error types (e.g., network, parsing, sandboxing).

- ■ `data_quality_score`: (For research data) A Guardian-computed score based on completeness, relevance, and format.
- ■ `bias_detection_flags`: (If applicable to generated code/data) Flags for detected biases.
- ■ `anomaly_flags`: (e.g., `HIGH_FAILURE_RATE`, `UNEXPECTED_RESOURCE_SPIKE`).
- ○ **Microscopic Detail:**
  - ■ **Threshold-Based Alerts:** Guardian generates `ToolQualityFeedbackSchema` messages when performance metrics for an Armory-provided tool cross predefined thresholds (e.g., success rate drops below 90%).
  - ■ **Armory's Adaptation:** This asynchronous feedback continuously informs Armory's `Dynamic Tool Registry & Discovery` and its `Learning Integration`. Armory might:
    - ■ Trigger a re-evaluation of the tool.
    - ■ Prioritize alternative tools for similar tasks.
    - ■ Flag the tool for maintenance or replacement.
    - ■ Refine its own internal heuristics for selecting the "best" tool for a given Core request.

### 7.1.5. Armory ↔ Script Kitty.Foundry: The Operational Substrate

Foundry provides the fundamental infrastructure services that enable Armory to function, from compute resources to data storage, deployment automation, and security. Armory is a significant consumer of Foundry's capabilities.

- ● **7.1.5.1. Armory to Foundry: Resource Provisioning via IaC Execution**

  - ○ **Purpose:** When Armory needs to provision external computational resources (VMs, specialized GPUs, quantum computing access, temporary storage buckets) for Script Kitty's operations (e.g., for a particularly demanding research task, or to spin up a custom environment for a rare tool), it leverages Foundry's `Resource Provisioning via IaC Templates` capability.
  - ○ **Communication Protocol:** Armory sends a `ResourceProvisioningRequestSchema` (Protobuf) to Foundry's dedicated gRPC endpoint for IaC execution.
  - ○ **Data Schema:** `ResourceProvisioningRequestSchema`.
    - ■ `request_id`, `task_id`: For correlation.
    - ■ `iac_template_id`: Reference to a pre-approved IaC template (e.g., "gpu-compute-instance", "s3-research-bucket").
    - ■ `parameters`: JSON object with template-specific variables (e.g., `region="us-east-1"`, `gpu_type="A100"`, `storage_size_gb=100`).

- **■** `expected_duration_hours`: Estimated lifespan of the resource.
- **■** `access_policy_id`: Reference to a pre-defined access policy for the provisioned resource.
- **■** `cost_estimate_threshold`: Max acceptable cost.
  - ○ **Microscopic Detail:**
    - ■ **Foundry's IaC Orchestration:** Foundry's `Kubernetes-Native Compute & Orchestration` layer receives this request. It uses internal `Terraform` or `Ansible` API clients to execute the specified IaC template, passing the provided parameters.
    - ■ **Credential Management:** Foundry retrieves necessary cloud provider credentials from its `Secrets Management` (HashiCorp Vault) for executing IaC.
    - ■ **Asynchronous Operation:** Resource provisioning is often long-running. Foundry provides `ResourceProvisioningStatusSchema` updates back to Armory via gRPC streams, informing Armory of progress (e.g., "VM creation in progress," "Resource provisioned, IP: 1.2.3.4").
    - ■ **Cost Monitoring (Guardian integration):** Foundry, during resource provisioning, integrates with Guardian's `Real-time Performance Monitoring` to track the actual cost of provisioned resources, flagging any exceeding thresholds.
- **●** **7.1.5.2. Armory leveraging Foundry's Data Management:**

  - ○ **Purpose:** Armory relies heavily on Foundry's `Data Management` sub-system for storing raw web data, processed insights, and generated tool artifacts.
  - ○ **Communication Protocol:** Standard S3 API calls (for object storage), DVC commands/API (for versioning), Feast API (for features).
  - ○ **Microscopic Detail:**
    - ■ **Object Storage:** Armory's `Intelligent Web Scraping & Data Digestion` components write raw HTML, extracted text, and parsed structured data directly to Foundry's `Scalable Object Storage` (e.g., MinIO or AWS S3). Files are organized by `source_url_hash`, `timestamp`, `data_type`.
    - ■ **DVC for Tool Code/Manifests:** When Armory generates a new tool wrapper or updates its `Dynamic Tool Registry` with a new `ToolManifestSchema`, it checks these into DVC-managed Git repositories, ensuring every version is tracked and reproducible. This links the generated code to its versioned metadata.
    - ■ **Feature Store for Tool Metadata:** Armory might publish metrics about tools (e.g., their average performance based on historical usage or their estimated resource consumption) as features into Foundry's `Feature`

`Store` (Feast). This allows Core or other agents to query these features for planning or decision-making.

- **7.1.5.3. Armory's Deployment and Runtime via Foundry's CI/CD & Kubernetes:**

  - **Purpose:** Foundry provides the foundational MLOps capabilities for building, deploying, and running Armory's microservices and LLM models.
  - **Communication Protocol:** Git commits trigger CI/CD pipelines; Kubernetes API for deployments.
  - **Microscopic Detail:**
    - **Dockerization:** Every Armory sub-component (e.g., Web Scraper service, Tool Wrapper Generation service, Tool Registry service) is containerized into a Docker image.
    - **CI/CD Pipelines:** Armory's codebase lives in Git repositories. Any commit triggers Foundry's `CI/CD Automation` (e.g., GitLab CI/CD), which:
      - Lints and tests Armory's code.
      - Builds new Docker images for Armory services.
      - Pushes these images to Foundry's `Secure Container Registry`.
      - Triggers the `CD/GitOps` process.
    - **GitOps Deployment:** Foundry's Argo CD/Flux CD instances monitor Armory's Git repository for Kubernetes manifests. When a new image tag is available, the GitOps agent automatically updates the Armory deployments in Kubernetes.
    - **Kubernetes Orchestration:** Foundry's `Kubernetes-Native Compute & Orchestration` manages Armory's pods:
      - **Scheduling:** Ensures Armory's web scraping jobs run on nodes with sufficient network access and resources.
      - **Scaling:** `Horizontal Pod Autoscalers` scale Armory's web scraping services based on queue depth of research requests or CPU utilization.
      - **Self-healing:** Kubernetes automatically restarts failed Armory containers/pods, maintaining service availability.
      - **Resource Allocation:** Enforces CPU/memory limits and requests for Armory pods, preventing resource contention.

- **7.1.5.4. Armory's Security and Observability via Foundry's Unified Stack:**

  - **Purpose:** Armory's operations are continuously monitored, secured, and audited by Foundry's comprehensive observability and security layers.
  - **Communication Protocol:** Internal logging/metrics APIs, network traffic for security policies.
  - **Microscopic Detail:**
    - **Centralized Logging:** Armory services emit structured logs (JSON format) to stdout/stderr. Foundry's `Fluent Bit` agents on each node

capture these logs and forward them to `Loki/Elasticsearch`. All Armory logs are enriched with `service_name="armory"`, `component_name="web_scraper"`, `request_id`, `task_id`.

- **Metrics:** Armory's services expose Prometheus metrics endpoints. Foundry's `Prometheus Operator` scrapes these metrics (e.g., `armory_scrape_success_total`, `armory_data_digestion_latency_seconds`, `armory_tool_wrapper_gen_errors_total`). These metrics are visualized in Grafana dashboards dedicated to Armory's performance.
- **Distributed Tracing:** Armory's code is instrumented with `OpenTelemetry SDKs`. Traces for each request (e.g., a web scraping job) are exported to Foundry's `Jaeger/Zipkin` backend, allowing developers to see the full execution path within Armory and its interactions with external APIs.
- **Secrets Management:** Armory services retrieve sensitive credentials (e.g., API keys for external services, cloud provider access keys for IaC execution) dynamically from Foundry's `HashiCorp Vault` via a secure client library, avoiding hardcoding secrets.
- **Network Policy:** Foundry's `Calico/Cilium` enforces `NetworkPolicy` rules, ensuring Armory's web scraping service is the only one allowed to initiate outbound connections to the public internet, and limiting internal communication to only authorized services (e.g., talking to Core, writing to Object Storage).
- **Runtime Security:** Foundry's `Falco` monitors Armory's containers for anomalous syscalls, such as attempts to access unauthorized files, spawn unexpected processes, or make suspicious network connections from the sandboxed environment.
- **Container Scanning:** New Docker images built for Armory's services are automatically scanned by Foundry's `Trivy/Clair` before deployment, ensuring no known vulnerabilities are introduced.

# Script Kitty.Skills

# Part IV: Script Kitty.Skills - The Task Execution AI

**Core Purpose:** Script Kitty.Skills represents the system's specialized effector arm, housing a diverse and dynamic collection of AI models and computational agents engineered for precise, high-fidelity execution of discrete tasks. It is the "doer" component, receiving granular instructions from Script Kitty.Core and transforming abstract plans into tangible results. This module is paramount for delivering on Script Kitty's promise of practical, real-world utility, encompassing everything from complex data analysis and robust code generation to advanced image processing and bespoke problem-solving. It is here that raw computational power meets targeted intelligence, ensuring every instruction from Core is executed with surgical precision and unwavering reliability, without losing any functionality previously outlined.

**Technical Stack Overview (Skills):**

- **Languages:** Python (dominant for ML models, data processing, and scripting), Go (for high-performance microservices, especially for sandboxed execution orchestration).
- **ML Frameworks:** PyTorch, TensorFlow, scikit-learn (for model training and inference).
- **Data Processing:** Pandas, NumPy, SciPy, Dask, Polars, DuckDB.
- **Model Serving:** TorchServe, TensorFlow Serving, KServe, Ray Serve, ONNX Runtime.
- **Sandboxing:** Docker, Linux Namespaces & cgroups, Firejail/Bubblewrap, seccomp-bpf.
- **Code Execution:** Jupyter Kernel Gateway, secure `subprocess` invocations.
- **Generative AI:** Specialized Code LLMs (e.g., CodeLlama, StarCoder2) served by vLLM/Text Generation Inference.
- **Testing:** pytest (for Python code validation).
- **Communication:** gRPC (for task receipt and result transmission), Kafka (for asynchronous events and large data streams).

---

**Minute Aspects & Technical Dissection:**

**1. Specialized AI Models & Inference Sub-system**

This sub-system is the operational hub for Script Kitty's diverse range of AI models, ensuring they are served efficiently, scalably, and securely to execute specific computational tasks identified by Script Kitty.Core.

- **1.1. Model Registry & Versioning (Integrated with Foundry):**

  - **Purpose:** A centralized, version-controlled repository for all trained models, metadata, and associated artifacts. Ensures reproducibility and manageability of models.
  - **Technical Dissection:**

- **Model Metadata:** Each registered model includes metadata such as `model_id` (UUID), `version`, `model_type` (e.g., `COMPUTER_VISION_CLASSIFICATION`, `NLP_SENTIMENT_ANALYSIS`, `TIME_SERIES_PREDICTION`), `input_schema` (Protobuf), `output_schema` (Protobuf), `resource_requirements` (CPU, GPU, memory), `dependencies` (Python packages, system libraries), `training_dataset_id`, `performance_metrics` (accuracy, latency), `last_updated_timestamp`, and `owner_agent` (e.g., "Skills.ImageProcessor").
            - **Storage:** Model weights and associated files are stored in object storage (e.g., MinIO or S3-compatible storage managed by Foundry) with strong consistency.
            - **Versioning:** Every model iteration (fine-tuning, re-training) receives a new, immutable version ID, allowing for rollbacks and A/B testing. Foundry's DVC (Data Version Control) integrates here to version model artifacts alongside data.
            - **API/SDK:** Provides a gRPC API for Skills agents to query available models, retrieve model paths, and access associated metadata.
- **1.2. Dynamic Model Loading & Inference Service (SkillExecutor Microservices):**

    - **Purpose:** Manages the lifecycle of model instances, dynamically loading them into memory and executing inference requests with optimal performance and resource utilization.
    - **Technical Dissection:**
        - **Containerized Serving:** Each specialized AI model or a logical grouping of models (e.g., all image classification models) is encapsulated within its own ephemeral Docker container. These containers are orchestrated by Kubernetes.
        - **Model Serving Frameworks:**
            - **TorchServe / TensorFlow Serving:** For deep learning models, these frameworks provide optimized HTTP/gRPC inference endpoints, batching capabilities, and model hot-reloading. They abstract away the underlying model graph/session management.
            - **KServe (KNative Serving for ML):** Leveraged for autoscaling (scaling to zero for inactive models, rapid scaling up on demand) and efficient serving of ML models within Kubernetes, integrating seamlessly with Istio for traffic management.
            - **Ray Serve:** For building complex, multi-model inference graphs or pipelines, and for serving Python-native models that benefit from Ray's distributed capabilities.
            - **ONNX Runtime:** Used for cross-framework model deployment, allowing models trained in PyTorch/TensorFlow to be converted to

ONNX format for optimized, portable inference across diverse hardware.

- **Resource Allocation:** Each SkillExecutor pod requests precise CPU, memory, and GPU resources from Kubernetes based on the `resource_requirements` specified in the model metadata. This prevents resource contention and ensures performance isolation.
- **Input/Output Validation:** Incoming inference requests are strictly validated against the model's `input_schema`. Output is similarly validated against the `output_schema` before being returned to Core.
- **Batching & Concurrency:** Serving frameworks are configured for optimal batching of incoming requests and concurrent inference executions to maximize GPU utilization.
- **Health Checks:** Kubernetes liveness and readiness probes continuously monitor the health of each SkillExecutor pod, ensuring only healthy instances receive traffic.
- **Error Handling:** Catches model-specific inference errors (e.g., out-of-memory, invalid input shape) and translates them into standardized `TaskResult` error messages for Core.

## 2. Sandboxed Code Execution & CLI Interaction Sub-system

This critical sub-system provides a secure, isolated, and controlled environment for executing arbitrary code or interacting with command-line interface (CLI) tools, preventing malicious actions or resource leaks.

- **2.1. Secure Execution Environment (CodeExecutor Microservice):**

  - **Purpose:** To create highly isolated, ephemeral containers for each code execution request, ensuring strong security boundaries and preventing privilege escalation or interference with other system components.
  - **Technical Dissection:**
    - **Ephemeral Docker Containers:** Every code execution request (e.g., a Python script generated by Skills, a CLI command) is launched within its own dedicated, short-lived Docker container. These containers are designed to be "throwaway" – created, executed, and immediately terminated.
    - **Linux Namespaces & cgroups:** Docker relies heavily on these kernel features for isolation:
      - **PID Namespace:** Isolates process IDs, so processes inside the container cannot see or affect processes outside.
      - **Network Namespace:** Provides a separate network stack, meaning the container has its own IP address and network interfaces, isolated from the host. Only explicitly allowed network calls (e.g., to internal APIs) are permitted.

- - - **Mount Namespace:** Isolates the filesystem, so the container has its own view of the filesystem hierarchy. Only necessary volumes (e.g., temporary storage for input/output files) are mounted.
    - **User Namespace:** Maps container UIDs/GIDs to different host UIDs/GIDs, enhancing isolation.
    - **cgroups (Control Groups):** Enforce strict resource limits on CPU, memory, I/O bandwidth, and process count for each container. This prevents a runaway script from consuming all system resources.
  - **Rootless Containers (Preference):** Where possible, containers are run as non-root users, even within the container, further reducing potential attack surface.
  - **Read-Only Root Filesystem:** The container's root filesystem is mounted as read-only, preventing persistence of malicious changes.
  - **Network Policies (Kubernetes):** Kubernetes NetworkPolicies, implemented by CNI plugins like Calico or Cilium, strictly define allowed ingress and egress traffic for `CodeExecutor` pods, preventing unauthorized external access or internal lateral movement.
  - **Seccomp-bpf Filters:** A custom `seccomp` profile is applied to each container. This Linux kernel feature filters system calls, allowing only a minimal whitelist of necessary system calls (e.g., `open`, `read`, `write`, `execve`). Any attempt to make an unapproved system call results in termination of the process, acting as a crucial last line of defense.
  - **AppArmor/SELinux Profiles:** Additional Mandatory Access Control (MAC) profiles (AppArmor for Ubuntu/Debian, SELinux for RHEL/CentOS) provide finer-grained control over what processes can do within the container (e.g., restrict file access, network access).
  - **Input/Output Redirection:** Standard input, output, and error streams from the executed code are redirected to capture logs and results.
  - **Timeout Mechanism:** A strict execution timeout is enforced for every task. If the code does not complete within the specified duration, the container is forcibly terminated.
- **2.2. Script Execution & CLI Orchestrator:**

  - **Purpose:** Manages the submission of code snippets and CLI commands to the sandboxed environment, retrieves results, and handles execution status.
  - **Technical Dissection:**
    - **Jupyter Kernel Gateway Integration:** For Python and R code execution, `Jupyter Kernel Gateway` can be leveraged. It provides an API to spin up isolated Jupyter kernels, submit code cells for execution, and retrieve outputs. Each kernel runs within its own sandboxed environment.
    - **subprocess Module (Python/Go):** For direct CLI command execution, the `subprocess` module in Python (or `os/exec` in Go) is used within the

`CodeExecutor` container. Inputs to CLI commands are carefully sanitized and validated to prevent command injection.

- **Argument Sanitization:** All command-line arguments and script inputs are rigorously sanitized and validated against expected formats and safe characters, preventing shell injection or path traversal attacks.
- **Result Capture:** Standard output, standard error, and exit codes of the executed process are captured and parsed. For Python scripts, structured output (e.g., JSON) is encouraged for easier parsing.
- **Error Reporting:** Distinguishes between execution errors (e.g., syntax errors, runtime exceptions) and sandboxing violations (e.g., timeout, security policy breach). Errors are meticulously reported back to Core via `TaskResult` messages.
- **File I/O Restriction:** Strictly limits file system access within the sandbox. Only explicitly granted temporary storage areas are writable. Data transfer (inputs, outputs) is typically done via in-memory buffers or secure, ephemeral volume mounts, rather than allowing arbitrary file system access.

## 3. Data Processing & Analysis Frameworks Sub-system

This sub-system provides the robust computational backbone for in-depth data manipulation, cleaning, statistical analysis, and visualization, operating on data provided by Core or Armory.

- **3.1. Data Transformation & Wrangling Agents:**

  - **Purpose:** Executes tasks requiring data cleaning, transformation, and preparation for analysis.
  - **Technical Dissection:**
    - **Python Ecosystem:** Leverages the full power of Python's scientific computing libraries:
      - **Pandas:** For tabular data manipulation (DataFrames), cleaning (handling missing values, duplicates), merging, filtering, and aggregation. Highly optimized C extensions for performance.
      - **NumPy:** For high-performance numerical operations on arrays and matrices, essential for vectorized computations.
      - **SciPy:** For advanced scientific and technical computing, including optimization, linear algebra, interpolation, and signal processing.
    - **Polars & DuckDB (for performance):** For tasks requiring extremely high performance on structured data, especially columnar operations. Polars (Rust-backed DataFrame library) and DuckDB (in-process analytical database) offer significant speedups over Pandas for certain workloads.
    - **Schema Enforcement:** Inputs and outputs of data transformation steps are validated against predefined data schemas (e.g., using Pydantic or Great Expectations) to ensure data quality and consistency.

- ■ **Distributed Processing (for Large Data):**
    - ■ **Dask:** For parallelizing Pandas, NumPy, and scikit-learn computations across multiple cores or nodes within a Kubernetes cluster, enabling analysis of datasets larger than single-machine memory. Dask workers are deployed as separate pods, coordinated by a Dask scheduler.
    - ■ **PySpark:** For extremely large datasets requiring full-fledged distributed processing, PySpark (Spark's Python API) can be integrated. Spark clusters can be dynamically provisioned by Foundry or run on existing infrastructure, with Skills agents submitting Spark jobs.
- ● **3.2. Statistical Analysis & Modeling Agents:**

    - ○ **Purpose:** Performs statistical inference, applies machine learning models for prediction or clustering, and generates insights.
    - ○ **Technical Dissection:**
        - ■ **Scikit-learn:** The workhorse for traditional machine learning algorithms (classification, regression, clustering, dimensionality reduction). Skills agents can be instructed to train new models (for simpler cases) or apply pre-trained models.
        - ■ **StatsModels:** For statistical modeling, hypothesis testing, and econometric analysis, providing robust statistical rigor.
        - ■ **Probabilistic Programming (Pyro, Edward2, Stan):** For tasks requiring Bayesian inference, uncertainty quantification, or complex generative models. These frameworks allow Script Kitty to reason about probabilities and integrate prior knowledge.
        - ■ **Time-Series Analysis:** Libraries like `Prophet` (Facebook) or `pmdarima` (auto-ARIMA) for forecasting and time-series specific analysis.
        - ■ **Feature Engineering:** The agents can be instructed to perform automated feature engineering using libraries like `featuretools` or custom transformation pipelines, often guided by Core's planning.
        - ■ **Model Evaluation:** Agents perform rigorous evaluation of trained models using appropriate metrics (e.g., accuracy, precision, recall, F1, RMSE, R-squared) and cross-validation techniques.
- ● **3.3. Data Visualization & Reporting Agents:**

    - ○ **Purpose:** Generates high-quality visualizations and structured reports based on data analysis results.
    - ○ **Technical Dissection:**
        - ■ **Matplotlib, Seaborn, Plotly:** Standard Python libraries for creating static and interactive plots.
        - ■ **Plotly Dash / Streamlit (for Interactive Dashboards):** If a dynamic, interactive visualization is required, Skills can generate a minimal Python

script using Dash or Streamlit, which Foundry then deploys as a temporary web application for Core to provide a link to the user.
- **Report Generation Libraries:** FPDF or ReportLab for generating PDF reports, or custom Markdown/HTML templating (e.g., Jinja2) for structured text reports.
- **Data Serialization for Reports:** Visualization data (e.g., plot configurations, aggregated tables) is serialized into standard formats (JSON, CSV, Parquet) for easy consumption by Nexus for display or by other agents.
- **Accessibility & Best Practices:** Adheres to data visualization best practices (e.g., clear labeling, appropriate chart types, colorblind-friendly palettes) to ensure clarity and interpretability of results.

## 4. Generative Task Execution (e.g., Programming, Report Generation) Sub-system

This sub-system leverages advanced generative AI models to create complex textual outputs, including code, reports, documentation, and other structured content, based on high-level instructions from Script Kitty.Core.

- **4.1. Code Generation & Refinement Agents:**

  - **Purpose:** Generates executable code in various programming languages (e.g., Python, SQL, shell scripts) based on natural language descriptions or functional specifications.
  - **Technical Dissection:**
    - **Code-Focused LLMs:** Utilizes specialized LLMs such as CodeLlama, StarCoder2, or DeepSeek Coder. These models are specifically pre-trained and fine-tuned on vast code corpora, making them highly proficient in syntax, common libraries, and programming paradigms.
    - **Contextual Code Generation:** The LLM receives detailed context from Core, including:
      - **Problem Description:** Natural language description of the desired functionality.
      - **Input/Output Specifications:** Expected data structures, APIs to use, or file formats.
      - **Constraints:** Performance requirements, specific algorithms to use, security considerations.
      - **Existing Code Snippets/Libraries:** Relevant parts of an existing codebase or available internal libraries to integrate with.
      - **Error Messages/Test Failures (for refinement):** If code generation is part of an iterative debugging loop, previous error messages or failed unit tests are fed back into the LLM as context for refinement.
    - **Static Analysis & Linting:** Generated code is immediately passed through static analysis tools (e.g., Bandit for Python security, ESLint for

JavaScript, SonarQube Community for general quality) and linters (e.g., Black, Flake8 for Python) to identify potential issues, security vulnerabilities, or style non-compliance.

- **Unit Test Generation & Execution:** The generative agent can also generate unit tests for the code it produced. These tests are then executed within the sandboxed environment (2.2) using frameworks like `pytest` (for Python). The test results (pass/fail, coverage) are fed back to the LLM for iterative refinement of the generated code. This forms a crucial self-correction loop.
- **Abstract Syntax Tree (AST) Parsing:** For critical code generation tasks, the generated code can be parsed into an AST (`ast` module in Python, `tree-sitter` for general parsing) to perform semantic validation and ensure structural correctness beyond simple syntax checks.
- **Human-in-the-Loop Integration (via Guardian):** For sensitive or highly complex code, the generated code and test results are forwarded to Guardian's policy enforcement engine, which may flag it for human review before final deployment or execution.

- **4.2. Structured Document & Report Generation Agents:**

  - **Purpose:** Generates comprehensive reports, summaries, technical documentation, or other structured textual content based on analytical results or research data.
  - **Technical Dissection:**
    - **Text Generation LLMs:** General-purpose LLMs (e.g., the same foundational LLM as Core, or specialized variants) are used, prompted to adhere to specific document structures (e.g., executive summary, methodology, results, conclusion).
    - **Content Integration:** Dynamically pulls data, visualizations, and insights from other Skills sub-systems (e.g., Data Analysis results, Armory's research findings).
    - **Templating and Markup:** Generates content in Markdown, LaTeX, HTML, or other formats suitable for easy rendering. Jinja2 templates can be used to define the overall structure of reports, with LLM-generated text populating specific sections.
    - **Tone & Style Control:** Prompts include directives on the desired tone (e.g., formal, informal, technical) and style, ensuring consistent output aligned with Script Kitty's persona or specific user requirements.
    - **Citation & Attribution:** If the report relies on external research or data, the agent is prompted to include proper citations and attribution.
    - **Versioning:** Generated documents are versioned and stored (e.g., in Foundry's artifact store) for auditability and future reference.

**Operational Excellence & Scalability (Skills-Specific):**

- **Elastic Scalability:**

  - **Microservice Architecture:** Each distinct Skill agent or logical grouping of models (e.g., all image classification, all code execution) is deployed as an independent microservice in Kubernetes.
  - **Horizontal Pod Autoscaling (HPA):** `SkillExecutor` and `CodeExecutor` pods are configured with HPAs, dynamically scaling up or down based on incoming request queue depth, CPU utilization, GPU utilization, or custom metrics (e.g., number of active code execution jobs).
  - **GPU Scheduling:** Kubernetes scheduling policies leverage node taints and tolerations, as well as resource quotas, to ensure GPU-intensive tasks are scheduled efficiently on GPU-enabled nodes.
  - **Dask/Spark Integration:** Dask and Spark clusters can be dynamically spun up by Foundry to provide on-demand distributed computing for large-scale data processing tasks, scaling worker nodes as needed.
- **Robustness & Resilience:**

  - **Health Checks:** Liveness and readiness probes are implemented for all Skills pods, ensuring unhealthy instances are removed from service and restarted.
  - **Automated Retries:** Task execution logic includes robust retry mechanisms with exponential backoff for transient failures (e.g., network issues, temporary resource unavailability).
  - **Circuit Breakers:** Implement circuit breakers in communication with Core to prevent cascading failures if a Skill service becomes overloaded or unresponsive.
  - **Idempotent Operations:** Design task execution to be largely idempotent where possible, allowing safe retries without unintended side effects.
- **Security & Auditability (Integrated with Guardian & Foundry):**

  - **Principle of Least Privilege:** Every Skill agent and sandboxed environment operates with the absolute minimum necessary permissions.
  - **Container Image Hardening:** Docker images for Skill services are built with minimal base images, few dependencies, and hardened security configurations.
  - **Secrets Management:** API keys or credentials required by Skills (e.g., for external APIs accessed by generated code) are securely managed and injected via HashiCorp Vault (or similar solution provided by Foundry), never hardcoded.
  - **Detailed Logging:** Comprehensive logging of all input requests, output results, execution times, resource utilization, and any errors/exceptions is channeled to Foundry's centralized logging system (Loki/Elasticsearch) for forensic analysis and auditing.
  - **Traceability:** OpenTelemetry tracing provides end-to-end visibility of every task executed by Skills, from its receipt from Core to its final result, including all internal sub-steps.

- **Continuous Improvement & Learning (Via Guardian & Foundry):**

  - **Feedback Loop:** Execution success/failure metrics, performance benchmarks, and any human corrections on generated code/reports (via Guardian's feedback interface) are continuously fed back to Script Kitty.Core.
  - **Automated Retraining Triggers:** Guardian monitors these metrics and can automatically trigger retraining pipelines (orchestrated by Foundry's Kubeflow Pipelines/Argo Workflows) for specific Skill models or generative LLMs if performance degradation is detected or new data becomes available.
  - **Dataset Versioning:** DVC (Data Version Control) ensures that all datasets used for training and evaluation within Skills are versioned and reproducible.
  - **Experiment Tracking:** MLflow is used to track all model training experiments, including hyperparameters, metrics, and model artifacts, allowing for effective model management and comparison.

# Skills interactions

**1. Script Kitty.Skills ↔ Script Kitty.Core: The Executor-Orchestrator Nexus**

This is the primary operational pipeline for task execution. Core delegates tasks, and Skills performs them, providing detailed updates.

- **1.1. Core to Skills: Task Delegation & Execution Request**

  - **Initiator:** Script Kitty.Core (specifically, its **Agent Orchestration & Lifecycle Management** sub-component, guided by the **Goal Decomposition & HTN Planning Engine**).
  - **Recipient:** Script Kitty.Skills (specifically, its **Specialized AI Models & Inference** or **Sandboxed Code Execution & CLI Interaction** sub-component).
  - **Communication Mechanism:** High-performance gRPC (Google Remote Procedure Call) via dedicated Protobuf messages. This ensures low-latency, schema-enforced, and bidirectional communication.
  - **Data Exchanged (`TaskExecutionRequestSchema` - Protobuf):**
    - `task_id` (string, UUID): A unique identifier generated by Core for tracking the entire lifecycle of this specific execution instance. Crucial for correlation across logs, metrics, and state.
    - `parent_plan_id` (string, UUID): Identifier of the overarching plan this task belongs to, enabling Skills to understand its context within Core's broader strategy.
    - `skill_type` (enum): Specifies the type of skill required (e.g., `CODE_EXECUTION`, `DATA_ANALYSIS`, `IMAGE_PROCESSING`, `LLM_GENERATION`, `MODEL_INFERENCE`, `QUANTUM_SIMULATION`). Guides Skills to route the request to the correct internal module.
    - `tool_identifier` (string): If the task requires a specific tool from Armory (e.g., "pandas_data_frame", "stable_diffusion_model"), its registered identifier.
    - `inputs` (map<string, AnyValue>): A map of dynamically typed parameters required by the skill. This is the heart of the execution payload. Examples:
      - For `CODE_EXECUTION`: `{"code_script": "import pandas...", "language": "python", "environment_config": {...}}`
      - For `MODEL_INFERENCE`: `{"model_name": "image_classifier_v2", "input_data": base64_encoded_image_bytes}`
      - For `DATA_ANALYSIS`: `{"dataset_uri": "s3://data-lake/raw/data.csv", "analysis_query": "calculate_average_sales(date_range)"}`

- - ■ `resource_requirements` (ResourceAllocationSchema): Detailed specification of compute needs (e.g., `cpu_cores`, `memory_gb`, `gpu_type`, `gpu_count`, `disk_gb`). This is generated by Core's planner and passed directly to Foundry for ephemeral resource provisioning.
      - ■ `timeout_seconds` (int): Maximum allowed execution time for the task. If exceeded, Skills will terminate the task and report a timeout error to Core.
      - ■ `callback_address` (string): The gRPC endpoint within Core where Skills should send updates and results.
      - ■ `security_context` (SecurityContextSchema): Contains information like the `user_id` who initiated the request (for audit logs and compliance), and any specific isolation requirements for the sandbox.
    - ○ **Purpose/Logic:** Core, having decomposed a high-level goal and identified a specific executable step, formally delegates this execution to Skills. The detailed `TaskExecutionRequestSchema` ensures that Skills receives all necessary context, inputs, and constraints to perform its function autonomously.
    - ○ **Microscopic Details:**
      - ■ **Core's Pre-computation:** Before sending, Core's planning engine might retrieve tool specifications from its **Global Knowledge Graph (GKGS)** (which syncs with Armory's tool registry) to populate `tool_identifier` and validate `inputs`. It consults its internal resource catalog (also part of GKGS) to determine optimal `resource_requirements`.
      - ■ **Skills' Initial Validation:** Upon receiving the request, Skills' `TaskExecutor` service performs immediate validation:
        - ■ Schema validation of the `TaskExecutionRequestSchema` against the Protobuf definition.
        - ■ Input validation: Checks for missing or malformed `inputs` based on the `skill_type` and `tool_identifier`.
        - ■ Security context verification: Ensures the request originates from an authorized Core instance.
      - ■ **Foundry Integration:** Skills does not provision resources directly. Instead, it informs Foundry (via a gRPC call to Foundry's **Kubernetes-Native Compute & Orchestration** sub-system or through a pre-allocated queue by Core) about the `resource_requirements` for the incoming `task_id`. This dynamic resource allocation is a core responsibility of Foundry.
      - ■ **Queuing:** If immediate execution is not possible (e.g., resource contention, queue full), Skills' `TaskExecutor` might place the task in an internal prioritized queue, sending an `ACKNOWLEDGED` status update to Core.
- ● **1.2. Skills to Core: Status Updates & Progress Reporting**

- ○ **Initiator:** Script Kitty.Skills (specifically, its internal `TaskExecutor` and the components running the actual task).
- ○ **Recipient:** Script Kitty.Core (specifically, its **Agent Orchestration & Lifecycle Management** sub-component).
- ○ **Communication Mechanism:** gRPC streaming or discrete gRPC calls, leveraging the `callback_address` provided in the initial request.
- ○ **Data Exchanged (`TaskUpdateSchema` - Protobuf):**
  - ■ `task_id` (string): The unique identifier of the task being updated.
  - ■ `status` (enum): Current state of the task (e.g., `QUEUED`, `PROVISIONING_RESOURCES`, `RUNNING`, `EXECUTING_SUBSTEP`, `HALTED_FOR_REVIEW`, `SUCCESS`, `FAILURE`, `TIMEOUT`, `CANCELLED`).
  - ■ `progress_percentage` (float): An optional, granular indicator of progress for long-running tasks (e.g., 0.0 to 100.0).
  - ■ `message` (string): A human-readable message detailing the current status or any interim observations (e.g., "Starting Python script execution," "Model inference 50% complete," "Waiting for sandbox environment setup"). This is vital for Nexus to provide user feedback.
  - ■ `metadata` (map<string, AnyValue>, optional): Key-value pairs for additional context (e.g., `current_iteration`, `estimated_time_remaining`).
  - ■ `log_snippet` (string, optional): A small excerpt from the task's console output or internal logs, useful for immediate debugging by Core or for display to a human reviewer.
- ○ **Purpose/Logic:** Core requires constant visibility into the execution status of delegated tasks to effectively manage the overall plan. These updates allow Core to adjust subsequent steps, trigger parallel tasks, or report progress back to Nexus.
- ○ **Microscopic Details:**
  - ■ **Granular Updates:** For tasks involving multiple stages (e.g., data loading, pre-processing, model training, evaluation), Skills sends distinct `EXECUTING_SUBSTEP` updates with specific messages.
  - ■ **Heartbeats:** For very long-running tasks, Skills may send periodic `RUNNING` updates to signal it's still active and prevent Core from timing out the task internally.
  - ■ **Logging:** Every `TaskUpdate` is simultaneously sent to Foundry's **Observability Stack** (Centralized Logging) for comprehensive audit trails and debugging.
- **1.3. Skills to Core: Task Completion/Failure Notification & Results**

  - ○ **Initiator:** Script Kitty.Skills (upon successful completion or definitive failure of a task).
  - ○ **Recipient:** Script Kitty.Core (Agent Orchestration & Lifecycle Management).

- ○ **Communication Mechanism:** gRPC call to the `callback_address`.
- ○ **Data Exchanged (TaskResultSchema - Protobuf):**
  - ■ `task_id` (string): Identifier of the completed task.
  - ■ `status` (enum): Final status, always `SUCCESS` or `FAILURE`.
  - ■ `output_data` (AnyValue, optional): The primary result of the task. Examples:
    - ■ For `CODE_EXECUTION`: `{"result_variable": 42.0, "output_file_uri": "s3://results/analysis.json"}`
    - ■ For `MODEL_INFERENCE`: `{"prediction": "cat", "confidence": 0.98}`
    - ■ For `LLM_GENERATION`: `{"generated_text": "Here's the Python script...", "token_count": 500}`
    - ■ For `DATA_ANALYSIS`: `{"summary_statistics": {"mean": 10.5, "std_dev": 2.1}, "visualization_uri": "s3://results/plot.png"}`
  - ■ `error_details` (ErrorSchema, optional): If `status` is `FAILURE`, detailed information about the error:
    - ■ `error_code` (enum): Categorized error (e.g., `EXECUTION_ERROR`, `TIMEOUT`, `INVALID_INPUT`, `RESOURCE_FAILURE`, `SECURITY_VIOLATION`).
    - ■ `error_message` (string): Human-readable error description.
    - ■ `stack_trace` (string, optional): Full stack trace for code execution errors.
    - ■ `diagnostic_logs_uri` (string, optional): URI to a larger log file in Foundry's data lake for deep debugging.
  - ■ `metrics` (map<string, float>, optional): Performance metrics for the task (e.g., `execution_time_seconds`, `gpu_memory_used_mb`, `cost_usd`). This is forwarded to Guardian for evaluation and training.
  - ■ `audit_log_uri` (string, optional): URI to the full, immutable audit log of the sandbox execution.
- ○ **Purpose/Logic:** This is the ultimate signal of task completion. Core uses these results to:
  - ■ Update the `parent_plan_id`'s state.
  - ■ Feed `output_data` as `inputs` to subsequent tasks in the HTN plan.
  - ■ Trigger conditional branching in the plan based on success/failure.
  - ■ Report overall plan progress back to Nexus.
  - ■ Forward metrics and audit logs to Guardian for performance evaluation and safety checks.
- ○ **Microscopic Details:**

- **Core's Result Processing:** Core's **Agent Orchestration** component parses the `TaskResultSchema`. If `SUCCESS`, it updates its internal representation of the plan's state and prepares the next task. If `FAILURE`, it triggers Core's **Feedback & Learning Integration** to analyze the failure (potentially triggering a replan or reporting to Guardian for human intervention).
- **Foundry Data Storage:** `output_data` (especially large files like processed datasets, generated images, or comprehensive reports) is stored in Foundry's **Scalable Object Storage** (Data Lake), and the URI is returned in the `output_data` field. This ensures efficient handling of large outputs.
- **Guardian Integration:** `metrics`, `error_details`, and `audit_log_uri` are immediately forwarded to Guardian's **Real-time Performance Monitoring & Anomaly Detection** and **Policy Enforcement Engine** for evaluation and security checks.

---

### 2. Script Kitty.Skills ↔ Script Kitty.Armory: Tool Provisioning & Execution Environment Setup

Skills relies on Armory for access to external tools and the definition of execution environments.

- **2.1. Skills to Armory: Tool/Environment Manifest Request (Indirect via Core)**

  - **Initiator:** Script Kitty.Skills (indirectly, by receiving a `TaskExecutionRequestSchema` from Core that specifies a `tool_identifier` or `environment_config`).
  - **Recipient:** Script Kitty.Armory (Tool Registry & Discovery).
  - **Communication Mechanism:** While Skills doesn't *directly* call Armory for a manifest during task execution (Core is responsible for this pre-computation as part of its planning), Skills *relies* on the accurate and up-to-date information that Core has pulled from Armory. If Skills encounters an unrecognized `tool_identifier`, it will report a `FAILURE` back to Core.
  - **Microscopic Details:**
    - **Core's Role:** Core's **Goal Decomposition & HTN Planning Engine** queries its **Global Knowledge Graph (GKGS)** (which itself is continuously synced with Armory's **Dynamic Tool Registry & Discovery** service) to retrieve the `ToolManifestSchema` for a given `tool_identifier`. This manifest contains the necessary `adapter_type`, `input_parameters`, `output_schema`, and `required_environment` that Core then includes in the `TaskExecutionRequestSchema` sent to Skills.

- **Skills' Interpretation:** Skills' `TaskExecutor` interprets the `required_environment` (e.g., "Python 3.9 with Pandas", "PyTorch 2.0 with CUDA 12") to select or provision the correct sandboxed environment from Foundry. If the `tool_identifier` refers to a specific ML model, Skills ensures that model is loaded and ready within its **Specialized AI Models & Inference** sub-component.
- **2.2. Armory to Skills: Tool/Resource Provisioning (Indirect via Foundry)**

  - **Initiator:** Script Kitty.Armory (after generating an IaC template or a specific resource configuration).
  - **Recipient:** Script Kitty.Skills (the actual runtime environment where the tool will operate).
  - **Communication Mechanism:** Indirectly through Foundry. Armory's **Resource Provisioning via IaC Templates** pushes IaC configurations to Foundry's **Kubernetes-Native Compute & Orchestration** and **Data Management** layers.
  - **Microscopic Details:**
    - **Pre-allocation/On-demand:** For sandboxed code execution, Armory might generate a Terraform plan for a specialized VM or a Kubernetes pod configuration (e.g., with specific GPU types, memory limits, network access rules) which Foundry then provisions. Skills then executes within this provisioned environment.
    - **Tool Wrapper Deployment:** If Armory generates a custom tool wrapper (e.g., a Python script to interact with a niche API), Foundry's CI/CD pipeline deploys this wrapper as a new microservice or container image within Skills' execution environment, making it available to Skills' `Sandboxed Code Execution` component. This involves updating Skills' internal `ToolRegistry` (a local cache of available executable tools).
    - **Data Access:** Armory provisions storage buckets (via Foundry's Data Lake) for Skills to read input data from and write output data to. The URI for this data is passed by Core in the `TaskExecutionRequestSchema`.

---

### 3. Script Kitty.Skills ↔ Script Kitty.Guardian: Safety, Evaluation & Feedback

Guardian is the conscience and quality control for Skills' executions.

- **3.1. Skills to Guardian: Runtime Telemetry & Execution Artifacts**

  - **Initiator:** Script Kitty.Skills (its **Sandboxed Code Execution**, **Specialized AI Models & Inference**, and **Data Processing** sub-components).

- **Recipient:** Script Kitty.Guardian (its **Real-time Performance Monitoring & Anomaly Detection** and **Policy Enforcement Engine**).
- **Communication Mechanism:** Asynchronous message queue (e.g., Apache Kafka/Pulsar via Foundry's **Internal Communication Hub & Message Broker**), for high-throughput, decoupled event streaming. Also direct gRPC for critical alerts.
- **Data Exchanged (`ExecutionTelemetrySchema` - Protobuf):**
    - `task_id` (string): The identifier of the execution.
    - `event_type` (enum): `RESOURCE_USAGE`, `LOG_STREAM`, `SECURITY_EVENT`, `OUTPUT_ARTIFACT_METADATA`.
    - `timestamp` (long): When the event occurred.
    - `cpu_utilization_percent` (float), `memory_used_mb` (float), `gpu_utilization_percent` (float): Granular resource usage metrics. Sent periodically.
    - `network_io_mb` (float): Network ingress/egress.
    - `disk_io_mb` (float): Disk read/write activity.
    - `log_level` (enum), `log_message` (string): Streams from the sandboxed environment's stdout/stderr.
    - `security_violation_type` (enum): `SANDBOX_ESCAPE_ATTEMPT`, `UNAUTHORIZED_FILE_ACCESS`, `PROHIBITED_NETWORK_CALL`. Triggered by low-level kernel hooks (like Falco integration).
    - `output_artifact_uri` (string): URI to output data in the data lake (e.g., `s3://skills-results/task-xyz/plot.png`).
    - `model_confidence` (float): For ML inference tasks.
    - `generative_model_output_hash` (string): A cryptographic hash of any generated text/code for integrity checks and de-duplication.
- **Purpose/Logic:** Guardian needs a complete picture of Skills' behavior for safety, performance, and evaluation. This stream of telemetry is the raw data for its analysis.
- **Microscopic Details:**
    - **Foundry's Role:** Foundry's **Observability Stack** (Prometheus, Loki, OpenTelemetry) directly collects these metrics and logs from the Skills execution environment, and then forwards them to Guardian's dedicated Kafka topic.
    - **Low-Level Monitoring:** Skills' sandboxed environment is instrumented with kernel-level security tools (like Falco, managed by Foundry's **Runtime Security**). Any detected anomaly (e.g., an attempted sandbox escape) triggers an immediate high-priority alert to Guardian's Policy Enforcement Engine.
    - **Output Content Analysis:** For generative tasks, Skills might hash the output (`generative_model_output_hash`) and the full output content is sent to Guardian for **Bias Detection & Mitigation** and ethical policy

checks (e.g., checking if the generated text is harmful, biased, or violates content guidelines). This happens after the task's final output is sent to Core but *before* Core fully accepts it as valid for subsequent steps.

- **3.2. Guardian to Skills: Policy Enforcement & Action Termination/Correction**

  - **Initiator:** Script Kitty.Guardian (its **Policy Enforcement Engine**).
  - **Recipient:** Script Kitty.Skills (the running `TaskExecutor` for a specific `task_id`).
  - **Communication Mechanism:** Direct gRPC call for immediate, critical actions.
  - **Data Exchanged (`EnforcementActionSchema` - Protobuf):**
    - `task_id` (string): Identifier of the task to act upon.
    - `action` (enum): `TERMINATE_EXECUTION`, `PAUSE_EXECUTION`, `REQUEST_REPLAN`, `REQUEST_HUMAN_REVIEW`, `FLAG_OUTPUT_FOR_REVIEW`, `ADJUST_PARAMETERS_AND_RETRY`.
    - `reason` (string): Explanation for the action (e.g., "Resource over-utilization detected," "Generated content violated safety policy," "Detected sandbox escape attempt").
    - `severity` (enum): `CRITICAL`, `WARNING`, `ADVISORY`.
    - `suggested_parameters` (map<string, AnyValue>, optional): For `ADJUST_PARAMETERS_AND_RETRY`, providing new parameters for the skill.
  - **Purpose/Logic:** Guardian actively monitors Skills' behavior and outputs. If a policy violation, security threat, or critical resource issue is detected, Guardian can immediately intervene to prevent harm or correct the execution.
  - **Microscopic Details:**
    - **Pre-execution Scan:** For generative tasks, Guardian's Policy Enforcement Engine may perform a pre-execution safety check on the `inputs` that Core provides to Skills, potentially returning an `ADVISORY` or `REJECT` before execution even begins.
    - **Real-time Intervention:** Upon receiving a `TERMINATE_EXECUTION` action, Skills' `TaskExecutor` immediately kills the underlying process/container running the task, cleans up resources, and reports a `FAILURE` to Core with the `error_details` indicating Guardian's intervention.
    - **Human-in-the-Loop:** If `REQUEST_HUMAN_REVIEW` is sent, Skills pauses the execution and reports the status to Core, which then routes it to Nexus for displaying to the human operator via Guardian's **Human Oversight & Feedback Interface**.
    - **Feedback Loop:** This interaction is critical for Guardian's **Automated Evaluation & Training Trigger**. Failed tasks due to Guardian intervention (or successful ones with high confidence) are fed into Guardian's learning loop to refine policies or trigger model retraining.

## 4. Script Kitty.Skills ↔ Script Kitty.Foundry: Infrastructure & Runtime Fabric

Foundry is the literal ground beneath Skills' feet, providing all necessary compute, storage, networking, and operational services.

- **4.1. Skills to Foundry: Resource Request & Configuration Fetch**

  - **Initiator:** Script Kitty.Skills (`TaskExecutor`).
  - **Recipient:** Script Kitty.Foundry (primarily its **Kubernetes-Native Compute & Orchestration** and **Secrets Management** sub-systems).
  - **Communication Mechanism:** Kubernetes API calls (indirectly via a pre-configured service account with appropriate RBAC), and gRPC/REST calls to Foundry's services (e.g., Vault).
  - **Data Exchanged:**
    - **Kubernetes API:** Skills' `TaskExecutor` interacts with the Kubernetes API to:
      - Request new Pods/Jobs for sandboxed execution based on `resource_requirements` from Core's `TaskExecutionRequestSchema`. This involves defining Pod manifests (`apiVersion`, `kind: Pod`, `metadata`, `spec` including `containers`, `resources`, `securityContext`, `volumes`, `nodeSelector`, `tolerations`).
      - Query for existing resource availability.
      - Delete completed or failed execution pods.
    - **Vault API (for Secrets):** Skills makes authenticated requests to HashiCorp Vault (Foundry's **Secrets Management**) to dynamically fetch credentials required for specific tasks (e.g., database connection strings, API keys for external services, cloud storage access tokens).
    - **Feast API (for Features):** Skills' **Specialized AI Models & Inference** components might query Feast (Foundry's **Feature Store**) for real-time features required by their models (e.g., `get_online_features(entity_key={'user_id': 'xyz'}, feature_names=['user_activity_score', 'recent_purchase_history']))`.
  - **Purpose/Logic:** Skills cannot function without computational resources or secure access to necessary credentials and data. Foundry provides this on-demand.
  - **Microscopic Details:**
    - **Ephemeral Pods:** For each `CODE_EXECUTION` or `DATA_ANALYSIS` task, Skills dynamically requests a new, ephemeral Kubernetes Pod from Foundry. This Pod is pre-configured with the required environment

(Python, R, specific libraries) via a base Docker image from Foundry's **Container Registry**.

- **Strict Security Context:** The Pod's `securityContext` is rigorously defined to enforce least privilege, including:
    - `runAsNonRoot: true`
    - `allowPrivilegeEscalation: false`
    - `capabilities: drop all` (and then selectively add only necessary ones, e.g., `NET_RAW` if required for network analysis, but heavily scrutinized by Guardian).
    - `seccompProfile`: A strict Linux `seccomp` profile (managed by Foundry's **Security & Identity Management**) to whitelist allowed system calls, preventing unauthorized kernel interactions.
    - `readOnlyRootFilesystem: true`
- **Volume Mounts:** Ephemeral volumes are mounted for temporary task data, and persistent volumes (from Foundry's **Data Management**) are mounted for accessing datasets from the Data Lake.
- **Network Policies:** Foundry's **Network Policy** (Calico/Cilium) ensures that these ephemeral Skills pods can only initiate connections to allowed internal services (e.g., Core, Guardian, Foundry's data services) and explicitly whitelisted external endpoints (e.g., a specific public API required by the task), effectively sandboxing network access.
- **Service Mesh Sidecar:** An Istio/Linkerd proxy (from Foundry's **Service Mesh**) is injected into each Skills execution pod, enforcing mTLS for all outbound/inbound communication and providing granular authorization policies.

- **4.2. Foundry to Skills: Resource Provisioning & Operational Services**

    - **Initiator:** Script Kitty.Foundry (Kubernetes scheduler, Kubelet, various Foundry services).
    - **Recipient:** Script Kitty.Skills (the newly provisioned Pods/Job, or the Skills microservice itself).
    - **Communication Mechanism:** Kubernetes's internal mechanisms (API server, Kubelet), service mesh, object storage.
    - **Data Exchanged:**
        - **Pod Assignment & Execution:** Kubernetes scheduler assigns the Skills execution Pod to a suitable node. The Kubelet then pulls the specified Docker image from Foundry's **Container Registry** and launches the container within the defined security context.
        - **Environment Variables & ConfigMaps:** Foundry injects necessary configuration (e.g., `TASK_ID`, `CORE_CALLBACK_URL`, `VAULT_ADDR`, `FEAST_ADDR`) as environment variables or mounts them as `ConfigMaps` into the Skills execution pod.

- **Secrets Delivery:** Vault agents or CSI (Container Storage Interface) Secret Store drivers (managed by Foundry) securely inject secrets as files or environment variables into the Skills pod *only* at runtime, ensuring they are never hardcoded.
- **Data Access URIs:** Foundry's **Data Management** components provide pre-signed URLs or S3-compatible endpoints for Skills to securely read/write data from/to the Data Lake.
- **Telemetry Agents:** Foundry ensures that Prometheus exporters, Fluent Bit/Fluentd agents, and OpenTelemetry SDKs are running within or alongside every Skills execution container, continuously streaming metrics, logs, and traces back to Foundry's **Observability Stack**.
- ○ **Purpose/Logic:** Foundry orchestrates the entire runtime environment for Skills, ensuring it has everything it needs to operate securely and efficiently.
- ○ **Microscopic Details:**
  - **Ephemeral Nature:** Foundry's design prioritizes the ephemeral nature of Skills' execution pods. Once a task is complete (success or failure), the pod is automatically terminated and garbage collected, minimizing resource waste and attack surface.
  - **Resource Quotas:** Foundry applies `ResourceQuotas` and `LimitRanges` to the namespaces where Skills tasks execute, preventing any single task from consuming excessive cluster resources.
  - **Service Mesh:** Foundry's **Service Mesh** (Istio/Linkerd) ensures that the Skills execution pod's network interactions are governed by mTLS and strict access policies, enforced at the proxy sidecar injected into the pod.

---

### 5. Script Kitty.Skills ↔ Script Kitty.Nexus: Indirect User Interaction

Skills does not directly interact with Nexus. All its outputs and progress updates are channeled through Core, which then synthesizes them for Nexus.

- **5.1. Skills to Core to Nexus: Task Progress & Results for User Display**
  - ○ **Initiator:** Script Kitty.Skills (via `TaskUpdateSchema` and `TaskResultSchema`).
  - ○ **Intermediate:** Script Kitty.Core (processes updates/results, updates plan state).
  - ○ **Recipient:** Script Kitty.Nexus (Response Synthesis & Persona Layer).
  - ○ **Communication Mechanism:** Skills -> Core (gRPC). Core -> Nexus (gRPC, standard `ScriptKittyMessageSchema`).
  - ○ **Microscopic Details:**
    - **Core's Synthesis:** Core's **Agent Orchestration** component receives raw `TaskUpdate` and `TaskResult` messages from Skills. Instead of blindly forwarding, Core synthesizes these into high-level, user-friendly `message` strings for Nexus. For example, multiple granular `EXECUTING_SUBSTEP`

updates from Skills might be consolidated into a single "Skills is performing data analysis..." message for Nexus.

- **Attribution:** Core explicitly includes `source_agent: "Skills"` in the `ScriptKittyMessageSchema` it sends to Nexus. This allows Nexus's **Group Chat Facilitation & Multi-Agent Attribution** to prepend "Skills reports:" or "Skills has completed the task:" before rendering the final text.
- **Error Handling:** If Skills reports a `FAILURE`, Core determines if this is a recoverable error (triggering a replan) or a critical failure requiring immediate user notification. It then crafts an appropriate error message for Nexus to display.
- **Result Formatting:** Core might also format the raw `output_data` from Skills into a more presentable format (e.g., generating Markdown tables or summary sentences) before sending it to Nexus for display to the user.

Script Kitty.Guardian

## Part 5: Script Kitty.Guardian - Safety & Alignment Proxy / Evaluation & Training AI

**Core Purpose:** Script Kitty.Guardian stands as the unwavering conscience and the relentless self-improvement engine of the entire Script Kitty ecosystem. It is meticulously designed to enforce a robust matrix of ethical, legal, and safety controls, continuously monitor the system's performance and behavior, proactively identify learning opportunities, and strategically trigger comprehensive training and retraining cycles across all Script Kitty models and modules. Guardian is the ultimate guarantor of Script Kitty's responsible evolution and its unwavering alignment with human values.

**Technical Stack Overview (Guardian):**

- **Language:** Python (for policy logic, ML pipelines, and data processing), Go (for high-performance policy enforcement and stream processing agents).
- **Policy Enforcement:** Open Policy Agent (OPA) for declarative policies, with custom Go/Python rule engines for complex scenarios.
- **Stream Processing:** Apache Flink (or Apache Spark Streaming) for real-time data ingestion and anomaly detection.
- **Databases:** PostgreSQL (for policy rules, audit logs, evaluation results, and long-term storage), Redis (for real-time metrics and temporary state).
- **MLOps Orchestration:** Kubeflow Pipelines or Argo Workflows for managing complex ML training and evaluation pipelines.
- **Experiment Tracking:** MLflow for logging model parameters, metrics, and artifacts.
- **Data Versioning:** DVC (Data Version Control) for reproducible datasets and models.
- **Bias Detection:** Aequitas, Fairlearn, IBM's AI Fairness 360 (AIF360).
- **Observability:** Prometheus, Grafana, OpenTelemetry (integrated with Foundry).

---

**Minute Aspects & Technical Dissection:**

**1. Policy Enforcement Engine Sub-system**

This sub-system is the critical gatekeeper, intercepting and rigorously evaluating every proposed action against predefined ethical, legal, and safety policies. It's the primary point where human-in-the-loop intervention can be enforced for critical actions.

- **1.1. Action Interception & Serialization (gRPC Interceptor):**

  - **Purpose:** To capture every single "action request" generated by Script Kitty.Core or any other agent *before* it is executed.
  - **Technical Dissection:**

- **gRPC Interceptors:** Every gRPC endpoint exposed by Core or other agents, particularly those that trigger external actions (e.g., API calls, code execution, infrastructure provisioning), is instrumented with a server-side gRPC interceptor. This interceptor transparently intercepts the incoming `ActionRequestSchema` (a Protobuf-defined message detailing the proposed action, its parameters, originating agent, and criticality level).
- **Schema Enforcement:** The `ActionRequestSchema` is strictly defined to include all necessary metadata for policy evaluation: `action_type` (e.g., `API_CALL`, `CODE_EXECUTION`, `RESOURCE_PROVISION`), `target` (e.g., URL, function name), `parameters` (structured key-value pairs), `originating_agent_id`, `user_context`, `criticality_level` (e.g., `LOW`, `MEDIUM`, `HIGH`, `CRITICAL`), and `justification_summary`.
- **Secure Channel:** Communication between agents and Guardian's Policy Enforcement Engine utilizes mTLS (mutual TLS) via a service mesh (Istio/Linkerd) to ensure authenticated and encrypted communication, preventing unauthorized tampering with action requests.

- **1.2. Policy Evaluation & Decisioning (OPA + Custom Rule Engine):**

  - **Purpose:** To apply complex, declarative, and context-sensitive policy rules to determine if a proposed action is permissible.
  - **Technical Dissection:**
    - **Open Policy Agent (OPA) Integration:** OPA microservices run alongside Guardian, acting as policy decision points. The intercepted `ActionRequestSchema` is serialized into JSON and sent to OPA as an input document.
    - **Rego Policies:** Policies are written in OPA's high-level declarative language, Rego. These policies define `allow` or `deny` rules based on the `action_type`, `target`, `parameters`, `user_context` (e.g., "user is admin"), `criticality_level`, and real-time context fetched from the Global Knowledge Graph (GKGS) or external systems (e.g., "deny if target URL is on a blacklist," "require human approval if `criticality_level` is CRITICAL and `action_type` is `RESOURCE_PROVISION`").
    - **Policy Bundles:** Policies are managed as "bundles" in a Git repository (GitOps principle) and automatically distributed to OPA instances, ensuring version control and auditability.
    - **Custom Rule Engine (for nuanced ethical dilemmas):** For policies requiring more nuanced reasoning, particularly in complex ethical dilemmas where Rego's declarative nature might be insufficient, a custom Python/Go rule engine is employed. This engine might incorporate:

- **Graph Traversal:** Traversing the GKGS to assess potential second-order effects of an action.
- **Heuristic-based Scoring:** Applying weighted heuristics to assign a "risk score" to an action.
- **Ethical LLM (Fine-tuned):** A smaller, specialized LLM, fine-tuned on ethical frameworks (e.g., utilitarianism, deontology) and common ethical dilemmas, is used *not* to make a definitive decision, but to provide an *analysis* of potential ethical conflicts within the proposed action, flagging them for OPA or human review. Its output (e.g., "Potential for privacy violation due to data access," "Risk of misuse if code is generated without sandbox") is fed as an input context to OPA or the human reviewer.
- **Decision Output:** The evaluation yields a clear `PERMIT`, `DENY`, or `REQUIRE_HUMAN_REVIEW` decision, along with an `explanation` for the decision.
- **Decision Caching:** For frequently occurring, low-risk actions, OPA's decision caching can be utilized to improve performance, while ensuring cache invalidation for updated policies.

- **1.3. Human-in-the-Loop Orchestration:**

  - **Purpose:** To facilitate seamless human oversight and intervention for actions flagged for review.
  - **Technical Dissection:**
    - **Review Queue:** Actions flagged `REQUIRE_HUMAN_REVIEW` are placed in a persistent queue (e.g., Kafka topic, PostgreSQL table) and simultaneously pushed to a dedicated human review interface (a web dashboard).
    - **Reviewer Interface:** A user-friendly web interface (built with React/Next.js and a FastAPI backend) presents the `ActionRequestSchema` details, OPA's explanation, and any ethical LLM analysis. Reviewers can `APPROVE`, `DENY`, or `REQUEST_MORE_INFO`.
    - **Authentication & Authorization:** The review interface is secured with robust authentication (e.g., Keycloak) and role-based access control (RBAC), ensuring only authorized personnel can review specific types of actions.
    - **Asynchronous Feedback:** Once a human decision is made, it is asynchronously sent back to Guardian, which then signals Script Kitty.Core to either proceed or abort the action.
    - **Decision Audit:** Every human decision, along with the reviewer's identity and timestamp, is meticulously logged for full auditability.
- **1.4. Resilience & Fault Tolerance (Critical Design):**

- ○ **Purpose:** To ensure the Guardian itself is robust and that safety policies are maintained even in degraded operational states.
- ○ **Technical Dissection:**
  - ■ **Fail-Open vs. Fail-Closed Modes:** This is a crucial configuration choice, highly dependent on the `criticality_level` of the actions being governed:
    - ■ **Fail-Closed (Default for CRITICAL/HIGH actions):** If Guardian's Policy Enforcement Engine becomes unreachable or experiences severe errors for critical actions (e.g., `RESOURCE_PROVISION` for production systems, `DATA_DELETION`), all pending actions are automatically *blocked*. This prioritizes safety and prevents potentially catastrophic failures. A warning is logged, and an alert is immediately fired.
    - ■ **Fail-Open (Configurable for LOW/MEDIUM actions):** For less critical actions (e.g., `INFORMATIONAL_QUERY`, `GENERAL_CHAT` responses), Guardian can be configured to temporarily *allow* actions to proceed if its policy engine is unreachable. However, every such bypass is logged with a prominent warning and triggers an immediate high-priority audit event for post-mortem review. This prevents a Guardian outage from completely halting Script Kitty's non-critical operations.
  - ■ **Redundancy & High Availability:** Guardian's microservices (OPA instances, custom rule engines, Human-in-the-Loop backend) are deployed with multiple replicas across different availability zones/nodes within Kubernetes, ensuring no single point of failure.
  - ■ **Circuit Breakers:** All calls *to* Guardian's Policy Engine from other Script Kitty components are wrapped with circuit breakers (e.g., via a service mesh's capabilities). If Guardian's latency spikes or error rates exceed thresholds, the circuit breaks, preventing cascading failures and allowing for fallback mechanisms (like fail-open).
  - ■ **Self-Healing Kubernetes:** Kubernetes's native self-healing capabilities (liveness and readiness probes) automatically restart or reschedule unhealthy Guardian pods.
  - ■ **Graceful Degradation:** In a severely degraded state, Guardian might temporarily switch to a simpler, more restrictive set of emergency policies (e.g., "deny all external API calls except essential ones") to maintain basic safety until full functionality is restored.

## 2. Real-time Performance Monitoring & Anomaly Detection Sub-system

This sub-system acts as Script Kitty's nervous system, continuously collecting and analyzing operational data to ensure optimal performance and rapidly detect any deviations.

- **2.1. Metrics Collection & Aggregation (Prometheus & OpenTelemetry):**

  - **Purpose:** To gather comprehensive performance and operational metrics from every corner of Script Kitty.
  - **Technical Dissection:**
    - **Prometheus Exporters:** All Script Kitty microservices (Nexus, Core, Armory, Skills, Guardian, Foundry components) are instrumented with Prometheus client libraries, exposing custom metrics (e.g., `service_latency_seconds`, `error_total`, `task_success_rate`, `llm_token_generation_speed`, `gpu_utilization_percent`, `policy_evaluation_time`).
    - **OpenTelemetry:** Standardized instrumentation via OpenTelemetry SDKs is deployed across all services. This generates distributed traces for end-to-end request flows, allowing visualization of latency bottlenecks and service dependencies (e.g., Jaeger/Zipkin for visualization).
    - **Prometheus Server:** Deployed in Kubernetes, it scrapes metrics from all exposed endpoints at regular intervals.
    - **Metrics Storage:** Prometheus's time-series database for short-term, high-granularity metrics. Long-term metrics can be sent to a dedicated Mimir/Thanos cluster for extended retention and global queryability.
- **2.2. Log Aggregation & Analysis (Loki/Elasticsearch & Fluent Bit):**

  - **Purpose:** To centralize and make searchable all operational logs for debugging, auditing, and analysis.
  - **Technical Dissection:**
    - **Fluent Bit:** Deployed as a DaemonSet on every Kubernetes node, Fluent Bit efficiently collects logs from all containerized applications (standard output/error streams) and forwards them to a centralized logging backend.
    - **Loki:** A horizontally scalable log aggregation system, optimized for semi-structured logs. It stores logs as streams and uses Grafana for querying and visualization, making it highly efficient for operational debugging.
    - **Elasticsearch/OpenSearch:** For highly structured logs (e.g., audit trails, security events) where full-text search and complex aggregations are required.
    - **Log Correlation:** Logs are enriched with tracing IDs (from OpenTelemetry), session IDs, and user IDs, enabling correlation across distributed services.
- **2.3. Anomaly Detection Engine (Flink/Spark Streaming & ML Models):**

  - **Purpose:** To automatically detect unusual patterns or deviations from normal behavior in real-time metrics and logs, indicating potential issues.
  - **Technical Dissection:**

- **Stream Processing Framework:** Apache Flink (or Spark Streaming) is used to ingest real-time metric streams (from Kafka/Pulsar topics where Prometheus pushes data) and log streams. Flink's low-latency, stateful processing capabilities are ideal for this.
- **Statistical Models:** For basic anomaly detection (e.g., sudden spikes in error rates, unexpected drops in `task_success_rate`), statistical methods like moving averages, standard deviation analysis, or Exponentially Weighted Moving Average (EWMA) are applied directly in Flink jobs.
- **Machine Learning Models:** For more complex anomalies, specialized ML models are deployed within Flink jobs:
  - **Isolation Forest:** For detecting outliers in multi-dimensional metric data (e.g., identifying unusual combinations of CPU usage, memory, and network I/O for a specific service).
  - **One-Class SVM:** To identify deviations from a learned "normal" operational profile.
  - **Time-Series Anomaly Detection:** Models like Prophet or ARIMA-based methods are trained on historical performance data to predict expected ranges and flag deviations.
  - **Unsupervised Learning:** For log analysis, techniques like clustering or topic modeling on log messages can identify novel error patterns or malicious activities.
- **Alerting Integration:** Detected anomalies trigger immediate alerts via Prometheus Alertmanager, which then dispatches notifications to on-call teams and Script Kitty.Guardian's Human Oversight Interface.
- **Feedback Loop:** Anomaly alerts are also fed back into the `Feedback & Learning Integration` sub-system to inform potential retraining or system adjustments.

## 3. Automated Evaluation & Training Trigger Sub-system

This sub-system transforms insights from monitoring and human feedback into actionable triggers for system-wide improvement, ensuring Script Kitty's continuous evolution.

- **3.1. Evaluation Pipeline Orchestrator (Kubeflow/Argo Workflows):**

  - **Purpose:** To define, execute, and manage complex ML evaluation pipelines for all Script Kitty models.
  - **Technical Dissection:**
    - **Pipeline Definition:** Evaluation pipelines are defined as directed acyclic graphs (DAGs) using Kubeflow Pipelines SDK or Argo Workflows YAML. Each step in the pipeline is a containerized component (e.g., `data_preprocessing_step`, `model_inference_step`, `metrics_calculation_step`, `bias_detection_step`).

- ■ **Automated Triggers:** Pipelines are triggered based on:
    - ■ **Scheduled Intervals:** Regular evaluations (e.g., daily, weekly) of core models (Nexus LLM, Core LLM, Skills models).
    - ■ **Human Feedback:** Specific feedback from Nexus or Guardian's Human Oversight Interface flagging performance issues or undesirable behavior.
    - ■ **Anomaly Detection:** Alerts from the Anomaly Detection Engine indicating performance degradation (e.g., persistent drop in NLU accuracy, increased task failure rate).
    - ■ **New Data Availability:** When significant new, labeled datasets become available (e.g., from human corrections, newly acquired knowledge).
- ■ **Metrics Calculation:** Pipelines calculate various performance metrics: NLU accuracy, task success rate, generation quality (using ROUGE, BLEU, human ratings), response latency, resource consumption.
- ■ **Performance Baselines:** Current model performance metrics are compared against predefined baselines. If performance falls below a threshold, it triggers a retraining recommendation.

- ● **3.2. Retraining & Fine-tuning Trigger:**

    - ○ **Purpose:** Based on evaluation results, automatically initiate the training or fine-tuning process for specific Script Kitty models.
    - ○ **Technical Dissection:**
        - ■ **Threshold-based Triggering:** If the calculated performance metrics for a specific model (e.g., Nexus LLM's NLU accuracy) fall below a defined threshold, or if a significant number of human corrections accumulate for a particular model, a retraining event is triggered.
        - ■ **Targeted Retraining:** Guardian determines which specific model or module requires retraining based on the evaluation scope. This avoids retraining the entire system unnecessarily.
        - ■ **Pipeline Invocation:** The trigger invokes a dedicated training pipeline (also managed by Kubeflow Pipelines/Argo Workflows) for the target model. This pipeline includes steps for data preparation, model training (e.g., fine-tuning with LoRA/QLoRA for LLMs), validation, and model registration.
        - ■ **Resource Allocation:** The training pipeline ensures that necessary computational resources (GPUs) are provisioned via Script Kitty.Foundry before training commences.

- ● **3.3. Experiment Tracking & Model Registry (MLflow & DVC):**

    - ○ **Purpose:** To maintain a comprehensive record of all model training runs, versions, and performance, ensuring reproducibility and facilitating model selection.
    - ○ **Technical Dissection:**

- **MLflow:** Used to log all parameters, metrics, and artifacts (e.g., model weights, configuration files) for every training run. This allows for easy comparison of different training experiments.
- **Model Registry:** MLflow's Model Registry (or a similar solution like Seldon Core's) manages model versions, stages (e.g., `Staging`, `Production`), and approval workflows.
- **DVC (Data Version Control):** Crucially, DVC versions all datasets used in training and evaluation. This ensures that a specific model version can always be reproduced with the exact data it was trained on, providing full reproducibility and auditability for "how" Script Kitty learned.
- **Artifact Storage:** Model artifacts (weights, vocabularies) are stored in object storage (MinIO, S3) and referenced in MLflow.

## 4. Bias Detection & Mitigation Sub-system

This specialized component proactively identifies and helps mitigate biases within Script Kitty's models and outputs, fostering fairness and ethical AI.

- **4.1. Automated Bias Auditing Pipelines:**

  - **Purpose:** To regularly analyze training data and model outputs for potential biases.
  - **Technical Dissection:**
    - **Pipeline Integration:** Bias detection tools are integrated as distinct steps within the automated evaluation pipelines (3.1).
    - **Bias Detection Libraries:**
      - **Aequitas:** For analyzing bias in binary classifiers, particularly useful for decisions made by Core or Guardian's internal classifiers. It quantifies disparities in metrics like false positive rate, false negative rate, and true positive rate across different protected groups (e.g., gender, race, age).
      - **Fairlearn:** Provides various mitigation algorithms (e.g., reweighing, exponentiated gradient reduction) that can be applied during or after training to reduce observed biases.
      - **IBM AI Fairness 360 (AIF360):** Offers a comprehensive suite of metrics for detecting unwanted bias in machine learning models and algorithms for bias mitigation.
      - **Google's What-If Tool:** While often used for human exploration, its backend can be integrated into pipelines for generating performance slices across different demographic attributes.
    - **Sensitive Attributes:** The system identifies and tags sensitive attributes in data, respecting privacy, to perform bias analysis (e.g., inferring gender from names, but primarily relying on explicitly provided or inferred group labels where ethically permissible and privacy-compliant).

- - **Generative Bias Analysis:** For LLMs (Nexus, Core, Skills), specialized techniques analyze generated text for gender, racial, or cultural stereotypes, as well as inappropriate content generation, by comparing outputs to predefined bias lexicons or through adversarial prompting.
    - **Bias Reporting:** Automated reports detailing detected biases are generated and made available through the Human Oversight Interface and audit logs.
- **4.2. Mitigation Strategy Integration & Human Review:**

  - **Purpose:** To provide tools and facilitate human intervention for addressing identified biases.
  - **Technical Dissection:**
    - **Mitigation Recommender:** Based on the type and severity of detected bias, Guardian can recommend specific mitigation strategies (e.g., "re-weigh training data," "apply post-processing calibration," "fine-tune with debiased dataset").
    - **Human Review & Override:** Identified biases are presented to human reviewers on the Human Oversight Interface, allowing them to confirm the bias, approve mitigation strategies, or provide corrective feedback.
    - **Bias Correction Data:** Human corrections on biased outputs or biased training data are meticulously captured and fed back into the `Feedback & Learning Integration` sub-system, forming a critical dataset for future debiasing efforts and model retraining.

## 5. Human Oversight & Feedback Interface Sub-system

This interface provides a clear, actionable dashboard for human users to monitor Script Kitty's performance, intervene in critical situations, and provide crucial feedback.

- **5.1. Centralized Dashboard (Grafana & Custom UI):**

  - **Purpose:** To provide a single pane of glass for monitoring Script Kitty's health, performance, and ethical compliance.
  - **Technical Dissection:**
    - **Grafana Dashboards:** Highly customized Grafana dashboards are used to visualize all metrics and logs from Prometheus and Loki/Elasticsearch. This includes:
      - System health (CPU, memory, network, GPU utilization).
      - Task success rates, latency, error rates across all agents.
      - LLM inference metrics (tokens/second, latency, model confidence).
      - Policy enforcement statistics (deny rate, human review rate).
      - Bias detection reports and trends.
      - Audit logs and critical security alerts.

- **Custom Web UI (React/Next.js & FastAPI):** For more interactive elements beyond Grafana's capabilities, a custom web application provides:
  - **Human Review Queue:** Dedicated interfaces for action approval/denial, bias review, and corrective feedback submission.
  - **Policy Management:** A UI for viewing existing Rego policies, their version history, and potentially proposing new policy rules (pending a GitOps workflow).
  - **Audit Log Explorer:** An advanced search and filter interface for all audit logs.
  - **Model Performance Reports:** Detailed, drill-down reports on model evaluation results and bias analyses.
- **5.2. Feedback & Correction Submission:**

  - **Purpose:** To enable human users to provide explicit approval, override decisions, or submit corrective feedback, directly influencing Script Kitty's learning and improvement.
  - **Technical Dissection:**
    - **Structured Feedback Forms:** The UI provides structured forms for specific types of feedback:
      - **Intent Correction:** "Nexus misidentified my intent."
      - **Bad Response:** "This response was unhelpful/incorrect/toxic." (Allows highlighting text and suggesting corrections).
      - **Task Failure Analysis:** "Task failed because X, not Y." (Allows correcting root cause analysis).
      - **Bias Flagging:** "I detected bias in this output."
    - **API Endpoints:** Dedicated REST API endpoints in Guardian's backend receive these structured feedback payloads.
    - **Feedback Processing:** Feedback is timestamped, associated with the original session and task, and stored in PostgreSQL. It is then processed by the `Feedback & Learning Integration` sub-system to trigger evaluation or retraining.
    - **Human-in-the-Loop Datalogging:** Every human interaction and feedback point is meticulously logged and serves as a vital source of ground truth for Script Kitty's continuous learning.

## 6. Feedback & Learning Integration Sub-system

This is the ultimate nexus where all observations, signals, and human input converge to drive Script Kitty's overarching capacity, functionality, and knowledge reinforcement.

- **6.1. Unified Data Ingestion Pipeline (Kafka/Pulsar & Flink):**

- ○ **Purpose:** To consolidate all raw experience data, performance metrics, human feedback, and success/failure signals into a centralized, real-time processing stream.
- ○ **Technical Dissection:**
  - ■ **Message Brokers:** Apache Kafka or Apache Pulsar serve as the high-throughput, fault-tolerant message brokers. All Script Kitty components publish relevant events (e.g., `TaskCompletedEvent`, `TaskFailedEvent`, `UserFeedbackEvent`, `PolicyViolationEvent`, `AnomalyDetectedEvent`) to specific topics.
  - ■ **Stream Processing:** Apache Flink (or Spark Streaming) consumes these event streams. Flink jobs perform:
    - ■ **Event Transformation:** Normalizing event schemas for consistency.
    - ■ **Feature Engineering:** Extracting relevant features from raw events (e.g., `time_to_completion`, `error_type`).
    - ■ **Data Aggregation:** Aggregating data for specific metrics or time windows.
    - ■ **Data Validation:** Ensuring incoming data quality.
- ● **6.2. Reinforcement Learning (RL) & Policy Refinement:**

  - ○ **Purpose:** To dynamically refine Script Kitty's planning strategies, task execution policies, and overall decision-making based on success and failure signals.
  - ○ **Technical Dissection:**
    - ■ **Experience Replay Buffer:** A central component (e.g., a distributed database or a dedicated service) stores sequences of actions, observations, and rewards/penalties from executed tasks (both successes and failures). This acts as a global "experience replay buffer."
    - ■ **Offline RL Training:** Periodically, or when sufficient new experience data accumulates, an offline reinforcement learning agent (e.g., using Ray RLlib with algorithms like PPO, SAC, or DDPG) trains on the collected experience.
    - ■ **Reward Function Engineering:** Crafting robust reward functions is paramount. Rewards are assigned based on:
      - ■ **Task Success/Failure:** Explicit positive/negative rewards.
      - ■ **User Satisfaction:** Implicit rewards derived from positive user feedback or task completion.
      - ■ **Resource Efficiency:** Penalties for excessive resource consumption or long execution times.
      - ■ **Policy Compliance:** Penalties for violating ethical/safety policies (from Guardian).
      - ■ **Human Corrections:** Specific penalties for instances where human override was required due to a Script Kitty error.

- - **Policy Update:** The RL agent learns improved decision-making policies (e.g., for Core's planning engine, Skills' execution strategies). These learned policies are then translated into either:
    - - **New LLM Fine-tuning Data:** Generating synthetic training examples based on improved policies.
      - **Updates to Symbolic Rules:** Modifying or adding rules in Core's HTN planner.
      - **Direct Policy Deployment:** For RL-specific components within Skills or Core.
  - **A/B Testing (Phased Rollout):** Newly learned policies or fine-tuned models are often rolled out gradually (Canary deployments, A/B tests) to monitor real-world performance before full deployment, minimizing risk.
- **6.3. Global Knowledge Graph (GKGS) Update & Refinement:**

  - **Purpose:** To continuously update and enrich Script Kitty's central knowledge base based on new information, learned facts, and system experience.
  - **Technical Dissection:**
    - **Automated Knowledge Extraction:** From successful task executions (e.g., data digested by Armory, facts extracted by Skills), new entities, relationships, and facts are automatically extracted and proposed for addition to the GKGS.
    - **Truth Verification:** Proposed additions undergo a verification process, potentially involving cross-referencing with other sources or human review (via Guardian's interface) for high-confidence updates.
    - **Semantic Consistency:** New knowledge is integrated ensuring semantic consistency with existing knowledge in the graph.
    - **Knowledge Graph Embedding Retraining:** As the GKGS grows, its embeddings (learned representations of entities and relationships) are periodically re-trained to ensure the graph's semantic search and reasoning capabilities remain optimal. This improves RAG accuracy for all agents.
    - **Rule/Heuristic Refinement:** Based on observed patterns and RL feedback, existing heuristics or rules within the GKGS (e.g., for task decomposition) are refined or new ones are added.

**Auditability & Compliance (Guardian's Unyielding Mandate):**

Guardian's design fundamentally embeds auditability and compliance at its core.

- **Meticulous Audit Logging:** Every decision (policy enforcement, anomaly detected, training triggered, human intervention) is logged with immutable timestamps, originating components, and justifications. These logs are stored in secure, tamper-evident storage (e.g., a WORM-compliant object storage or a blockchain-like ledger for critical events), providing a complete, verifiable trail of Script Kitty's operations.

- **Clear Policy Definitions (Rego):** The use of Rego for policy definition offers unparalleled transparency. Policies are human-readable, version-controlled, and can be independently audited, making it clear *why* a particular action was permitted or denied.
- **Traceability:** The integration of OpenTelemetry ensures end-to-end traceability of every request and every decision, allowing auditors to reconstruct the full context of any action.
- **Responsible AI Practices:** By consistently monitoring for bias, enforcing ethical guidelines, and providing a human-in-the-loop for critical decisions, Guardian is the cornerstone for demonstrating Script Kitty's adherence to emerging AI regulations (e.g., EU AI Act, NIST AI Risk Management Framework) and fostering public trust.

**Scalability (Designed for Unprecedented Growth):**

Guardian is engineered for massive scale, essential for an AGI that continuously learns and expands its operational scope.

- **Stateless Policy Engine (OPA):** OPA instances are inherently stateless. Each evaluation request is independent. This allows for horizontal scaling of Guardian's policy services by simply adding more OPA pods in Kubernetes, handling an ever-increasing volume of policy checks without performance degradation. Load balancing across these instances is handled by the service mesh.
- **Stream Processing (Apache Flink/Spark Streaming):** Flink and Spark Streaming are built for distributed, high-throughput, low-latency stream processing. They scale horizontally by distributing processing across a cluster of nodes, making them inherently capable of handling massive real-time data streams from all Script Kitty components for monitoring and anomaly detection.
- **Distributed Databases:** PostgreSQL for audit logs and persistent state is configured for high availability (replication, failover clusters) and horizontal scalability (sharding, partitioning for large datasets). Redis for real-time caches and temporary data is deployed as a cluster, distributing data and load across multiple nodes. This ensures that data storage does not become a bottleneck as Script Kitty grows.
- **Kubernetes Horizontal Pod Autoscaler (HPA):** All Guardian microservices (Policy Engine, Anomaly Detection, Evaluation Orchestrator, Feedback Processors) are configured for HPA. They automatically scale their pod replicas up or down based on real-time metrics such as CPU utilization, memory consumption, message queue depth (for Kafka/Pulsar consumers), or custom metrics (e.g., `pending_policy_evaluations_count`), ensuring Guardian can adapt dynamically to fluctuating workloads.

**Continuous Improvement Loop (The Heartbeat of Script Kitty's Evolution):**

Guardian is, in essence, the very engine of Script Kitty's self-improvement. It embodies the core principle that Script Kitty is not a static artifact but an ever-evolving intelligence.

- **Closed-Loop System:** By continuously monitoring Script Kitty's internal operations and external interactions, evaluating its performance against defined objectives, identifying

areas of weakness or bias, and then directly triggering targeted retraining and mitigation efforts, Guardian creates a powerful, closed-loop feedback system.

- **Knowledge Reinforcement:** Every successful task execution, every human correction, every detected anomaly, and every learned policy is meticulously captured and fed back into the system's "memory" (GKGS and RL experience replay). This ensures that Script Kitty's "overarching capacity, functionality, knowledge reinforcement" is not merely a theoretical concept but a tangible, data-driven process.
- **Alignment Through Iteration:** This continuous improvement loop is fundamental to ensuring Script Kitty's progressive alignment with human values. As biases are detected and mitigated, as ethical policies are refined, and as the system learns from its mistakes, it gradually becomes safer, fairer, and more trustworthy over time. It's a pragmatic approach to "evolving" towards an increasingly benevolent and capable AGI.
- **Driving Adaptability:** This constant feedback and retraining mechanism enables Script Kitty to adapt to new domains, new user requirements, and new operational challenges without requiring manual re-engineering. It is the core mechanism by which Script Kitty truly "evolves."

# Guardian Interaction

**1. Guardian & Script Kitty.Nexus: The Conscience at the Gateway**

**Purpose of Interaction:** Guardian provides the initial layer of ethical, safety, and policy enforcement directly at the user-facing interface. It also receives early-stage feedback on user satisfaction and potential safety flags raised by Nexus, allowing for rapid human intervention or policy refinement.

**Interaction Flow: Nexus → Guardian (Pre-NLU Safety Check)**

- **Trigger:** User input is received by Nexus's `User Input & Pre-processing` sub-system.
- **Data Flow & Schema:**
    - **Message Type:** `UserQueryMessage` (Protobuf).
    - **Content:** `session_id` (string), `user_id` (string), `raw_text` (string), `channel_type` (enum), `timestamp` (int64).
    - **Nexus Action:** After `Input Normalization & Security Filter` (initial PII redaction, heuristic toxicity scoring), Nexus's `Core NLU & Intent Routing` component (specifically, the pre-NLU safety check path) *does not directly send* the full message to Guardian for *every* interaction to prevent latency. Instead, it might:
        - **Threshold-based Alerting:** If Nexus's *internal heuristic toxicity score* (from 1.2 in Nexus) exceeds a configured threshold, Nexus immediately sends a `PotentialSafetyViolation` event to Guardian's `Policy Enforcement Engine`.
        - **Human Override Request:** If Nexus detects a highly ambiguous or potentially problematic query that it cannot confidently process or route, it can submit a `HumanReviewRequest` to Guardian.
- **Communication Protocol:** Asynchronous Event Stream via Apache Kafka.
    - Nexus's `Input Normalization & Security Filter` (or its NLU component) publishes `PotentialSafetyViolation` or `HumanReviewRequest` events to a dedicated Kafka topic (e.g., `sk.guardian.safety.alerts`).
- **Guardian Processing (Policy Enforcement Engine):**
    - **Kafka Consumer:** Guardian's `Policy Enforcement Engine` (PEE) constantly consumes messages from `sk.guardian.safety.alerts`.
    - **Rule Evaluation:** Upon receiving a `PotentialSafetyViolation` event, the PEE immediately evaluates the `raw_text` against its pre-defined ethical and safety policies (Rego rules for OPA, custom expert system rules for nuanced cases).
    - **Action Determination:**

- If the text is deemed severely problematic (e.g., hate speech, illegal content), Guardian issues a `BlockResponse` back to Nexus.
- If it requires human review, it generates a `HumanReviewTask` (Protobuf) and sends it to the `Human Oversight & Feedback Interface`.
  - **Logging & Metrics:** All safety violations, policy evaluations, and actions taken are logged (to Foundry's centralized logging) and trigger metrics (e.g., `guardian_policy_violations_total`, `guardian_blocked_responses_total`).
- **Nexus Action (Response to Guardian):**
  - **Kafka Consumer:** Nexus also consumes from a Kafka topic (e.g., `sk.nexus.policy.actions`) where Guardian publishes its decisions.
  - **Blocking:** If Nexus receives a `BlockResponse` for a `session_id`, it immediately generates a pre-defined safe response (e.g., "I cannot fulfill this request as it violates my safety guidelines.") and dispatches it to the user via its `User Output Renderer & Channel Dispatcher`. The original query is discarded.
  - **Internal State Update:** Nexus updates the `session_context` to reflect the safety violation, potentially reducing the `trust_score` for the user or triggering further monitoring.

**Interaction Flow: Guardian → Nexus (Feedback/Alignment Updates)**

- **Trigger:** Guardian's `Automated Evaluation & Training Trigger` detects that a model or system behavior needs adjustment based on evaluation results or human feedback.
- **Data Flow & Schema:**
  - **Message Type:** `PersonaAdjustmentRequest` (Protobuf).
  - **Content:** `adjustment_type` (enum, e.g., `TONE`, `VERBOSITY`), `parameter_name` (string), `new_value` (string/float), `model_id` (string, target Nexus LLM).
  - **Guardian Action:** After analyzing evaluation data or processing human feedback, Guardian might issue a subtle `PersonaAdjustmentRequest` to fine-tune Nexus's conversational style. For example, if feedback indicates Nexus is too verbose.
- **Communication Protocol:** gRPC or internal API call if they reside in the same trust boundary, or Kafka for asynchronous updates. For model weights, Foundry's CI/CD would be used, but for runtime persona adjustments, a direct API might be used.
- **Nexus Processing (Response Synthesis & Persona Layer):**
  - **API Endpoint:** Nexus's `Response Synthesis & Persona Layer` exposes a secure internal gRPC endpoint (e.g., `/persona/adjust`).
  - **LLM Prompt Reconfiguration:** Nexus dynamically updates its system prompt parameters for the Nexus LLM based on the `PersonaAdjustmentRequest`.

This allows for subtle, real-time adjustments to its conversational output without full model retraining.
  - ○ **Logging:** Logs the persona adjustment for auditability.

**Interaction Flow: Nexus → Guardian (User Feedback Loop)**

- **Trigger:** Nexus presents a response to the user. User explicitly provides feedback (e.g., "thumbs up/down" button, explicit "feedback" command).
- **Data Flow & Schema:**
  - ○ **Message Type:** `UserFeedbackEvent` (Protobuf).
  - ○ **Content:** `session_id`, `user_id`, `response_id` (ID of the response generated by Nexus), `feedback_type` (enum: `POSITIVE`, `NEGATIVE`, `NEEDS_IMPROVEMENT`), `comment` (optional string), `timestamp`.
- **Nexus Action:** Nexus's `User Output Renderer & Channel Dispatcher` (or a dedicated feedback listener) captures explicit user feedback and publishes it.
- **Communication Protocol:** Asynchronous Event Stream via Apache Kafka.
  - ○ Published to a Kafka topic (e.g., `sk.guardian.user.feedback`).
- **Guardian Processing (Human Oversight & Feedback Interface / Automated Evaluation):**
  - ○ **Kafka Consumer:** Guardian's `Human Oversight & Feedback Interface` (HOFI) and `Automated Evaluation & Training Trigger` (AETT) consume from `sk.guardian.user.feedback`.
  - ○ **HOFI:** Displays feedback in a dashboard for human review and annotation.
  - ○ **AETT:** Aggregates feedback, identifies patterns (e.g., consistent negative feedback on a specific type of response from a Skills agent), and uses this data to trigger targeted evaluation pipelines or flag potential bias for `Bias Detection & Mitigation`.
  - ○ **GKGS Update:** Negative feedback, if sufficiently verified, can update confidence scores or add negative examples to the Global Knowledge Graph (GKGS) associated with specific planning heuristics or agent behaviors.

**Security Considerations for Nexus-Guardian Interactions:**

- **mTLS (Service Mesh):** All gRPC and Kafka communications between Nexus and Guardian are secured with mutual TLS (mTLS) enforced by Foundry's Service Mesh (Istio/Linkerd), ensuring encrypted and authenticated communication.
- **Network Policies:** Strict Kubernetes network policies (Calico/Cilium) ensure that Nexus pods can only connect to Guardian's Kafka topics and specific gRPC endpoints, and vice-versa, preventing unauthorized access.
- **Rate Limiting:** Guardian's Kafka consumers implement rate limiting to prevent being overwhelmed by a flood of safety alerts from Nexus.
- **Input Validation:** Both sides rigorously validate incoming Protobuf messages against their schema to prevent malformed data leading to system crashes or vulnerabilities.

- **Audit Logging:** Every interaction is meticulously logged, including timestamps, source, destination, and the nature of the message, for complete audit trails.

---

**2. Guardian & Script Kitty.Core: The Conscience and the Brain**

**Purpose of Interaction:** Guardian exerts strategic control over Core's planning and orchestration, ensuring alignment with safety policies and ethical guidelines. It also monitors Core's performance and learning loops, providing critical feedback for high-level system improvements.

**Interaction Flow: Core → Guardian (Action Proposal & Policy Check)**

- **Trigger:** After `Goal Decomposition & HTN Planning Engine` in Core generates a high-level `PlanSchema` (representing a sequence of actions involving other agents), or when Core is about to perform a critical internal decision (e.g., modifying its own planning heuristics).
- **Data Flow & Schema:**
    - **Message Type:** `ActionProposal` (Protobuf).
    - **Content:** `task_id` (string), `plan_id` (string), `proposed_actions` (list of `ActionDescriptor` Protobufs, each describing a sub-action: `agent_name`, `action_type`, `parameters`, `resource_requests`), `context_summary` (string, brief summary of the user's intent and current session context), `estimated_risk` (float, Core's self-assessment of risk).
- **Core Action:** Core's `Agent Orchestration & Lifecycle Management` component is mandated to submit any significant `ActionProposal` (or `DecisionProposal` for internal changes) to Guardian *before* execution.
- **Communication Protocol:** Synchronous gRPC.
    - Core's `Agent Orchestration & Lifecycle Management` component makes a blocking gRPC call to Guardian's `/policy_enforce` endpoint. Core waits for Guardian's verdict.
- **Guardian Processing (Policy Enforcement Engine):**
    - **gRPC Server:** Guardian's `Policy Enforcement Engine` (PEE) runs a gRPC server listening on `/policy_enforce`.
    - **OPA Evaluation:** The `proposed_actions` are passed to Open Policy Agent (OPA) for evaluation against Rego policies. These policies encode rules like:
        - "Cannot access external APIs without explicit user consent."
        - "Cannot generate content that advocates violence."
        - "Cannot provision resources exceeding a certain cost threshold without human approval."
        - "Cannot execute code on the host machine; must be sandboxed."

- - - "Cannot modify system-critical configurations without multi-factor approval."
    - **Ethical LLM Check:** For nuanced ethical dilemmas or ambiguous proposals, the PEE might also query its internal `Ethical LLM (Fine-tuned)` (a smaller, specialized LLM) with the `proposed_actions` and `context_summary`. The LLM's role is to flag potential ethical conflicts or biases, *not* to make the final decision, but to inform the OPA rules or human review process.
    - **Verdict:** Based on OPA rules and ethical LLM flags, Guardian returns a `PolicyVerdict` (Protobuf) message.
      - `verdict_type` (enum: `ALLOW`, `DENY`, `REQUEST_HUMAN_REVIEW`).
      - `reason` (string, e.g., "Violates PII policy," "Requires human approval for high cost").
      - `required_approvers` (list of strings, for human review).
    - **Logging & Metrics:** Each proposal and verdict is meticulously logged and metrics are incremented (`guardian_action_proposals_total`, `guardian_denials_total`, `guardian_human_review_requests_total`).
- **Core Action (Response to Guardian's Verdict):**
    - **Decision Branching:** Core receives the `PolicyVerdict`.
    - If `ALLOW`: Core proceeds to orchestrate the actions defined in the `PlanSchema`.
    - If `DENY`: Core immediately aborts the current plan, logs the denial reason, and may generate a new plan that adheres to the policies or informs Nexus (and thus the user) of the limitation.
    - If `REQUEST_HUMAN_REVIEW`: Core pauses the execution of the plan, updates its internal task state to `PENDING_HUMAN_REVIEW`, and informs Nexus to prompt the user or a human operator (via Guardian's HOFI) for approval.
- **Timeout & Retries:** Core implements timeouts for Guardian's response. If Guardian doesn't respond within a configured time, Core might retry or revert to a safe fallback plan, logging the incident.

**Interaction Flow: Guardian → Core (Direct Control/Updates - Policy/Learning)**

- **Trigger:** Guardian's `Automated Evaluation & Training Trigger` identifies a need to update Core's internal planning heuristics, model weights, or system configurations due to performance degradation, bias detection, or new learning.
- **Data Flow & Schema:**
    - **Message Type:** `SystemConfigUpdate` (Protobuf) or `ModelUpdateNotification` (Protobuf).
    - **Content (SystemConfigUpdate):** `target_component` (enum: `CORE`), `config_key` (string, e.g., `planning_heuristic_weights`), `new_value` (JSON/binary blob), `reason` (string).

- ○ **Content (ModelUpdateNotification):** `model_id` (string), `model_type` (enum, e.g., `CORE_PLANNING_LLM`), `new_image_tag` (string, pointing to a new Docker image in Foundry's registry), `reason`.
- **Guardian Action:**
  - ○ **Config Update:** Guardian can directly push `SystemConfigUpdate` messages to Core for immediate runtime adjustments of non-critical parameters (e.g., dynamically adjusting a confidence threshold for plan selection). This uses gRPC for a direct, authenticated call.
  - ○ **Model Retraining/Deployment Trigger:** For changes requiring a new model artifact (e.g., a fine-tuned Core planning LLM), Guardian doesn't push the model directly. Instead, it triggers a CI/CD pipeline in **Foundry** by submitting a `ModelRetrainRequest` to Foundry's `CI/CD & GitOps Automation` component (often via a Kafka event or a Foundry-exposed API). Foundry then builds and deploys the new model image, and subsequently notifies Core.
- **Core Processing (Feedback & Learning Integration / Internal Communication Hub):**
  - ○ **gRPC Endpoint:** Core's `Feedback & Learning Integration` component exposes a secure gRPC endpoint (`/system_config_update`) to receive `SystemConfigUpdate` messages.
  - ○ **Dynamic Reconfiguration:** Core's internal configuration store is updated based on the `config_key` and `new_value`. This triggers internal reconfiguration mechanisms.
  - ○ **Model Reload/Restart:** If a `ModelUpdateNotification` is received (from Foundry, after a new model is deployed), Core's `Internal Communication Hub` might orchestrate a graceful restart or hot-reload of its internal LLM inference service to pick up the new model image. This is coordinated with Foundry's Kubernetes orchestration.
  - ○ **Logging:** All updates are logged for audit and debugging.

**Interaction Flow: Core → Guardian (Performance Metrics & Audit Logs)**

- **Trigger:** Core continuously operates and executes tasks.
- **Data Flow & Schema:**
  - ○ **Metrics:** Core's `Agent Orchestration & Lifecycle Management` and `Goal Decomposition` components expose Prometheus metrics (e.g., `core_planning_latency_seconds`, `core_plan_success_rate`, `htn_decomposition_time_ms`, `core_llm_token_usage_total`).
  - ○ **Logs:** Core generates detailed structured logs (e.g., `plan_generated_event`, `agent_task_assigned_event`, `orchestration_error_event`, `core_llm_inference_details`) for every significant internal action.
  - ○ **Traces:** OpenTelemetry traces capture the flow of a task from planning through sub-agent orchestration.
- **Core Action:**

- **Metrics:** Prometheus (part of Foundry) scrapes metrics from Core's `/metrics` endpoint.
    - **Logs:** Core logs to stdout/stderr, which are collected by Fluent Bit/Fluentd (part of Foundry) and sent to Loki/Elasticsearch.
    - **Traces:** OpenTelemetry SDKs in Core export spans to the OpenTelemetry Collector (part of Foundry).
- **Guardian Processing (Real-time Performance Monitoring & Anomaly Detection / Automated Evaluation):**
    - **Prometheus/Grafana Integration:** Guardian accesses Core's metrics via Prometheus for `Real-time Performance Monitoring & Anomaly Detection` (RPMAD). Pre-configured Grafana dashboards display Core's health, and Prometheus Alertmanager triggers alerts if Core's performance deviates (e.g., planning latency spikes, success rate drops).
    - **Log Analysis:** RPMAD also processes Core's logs (from Loki/Elasticsearch) to identify specific error patterns or frequent failures that might indicate a systemic issue or a bias in Core's planning.
    - **Trace Analysis:** For complex issues flagged by metrics/logs, Guardian's team or automated systems can use Jaeger/Zipkin to inspect Core's traces, pinpointing the exact bottleneck or failure point in Core's orchestration.
    - **AETT Trigger:** RPMAD uses anomaly detection algorithms (e.g., Isolation Forest on `core_plan_success_rate` data streams) to identify degradation. If a significant anomaly is detected, AETT triggers a deeper evaluation pipeline, potentially leading to a Core retraining request sent to Foundry.
    - **Feedback Loop:** This continuous monitoring forms the primary feedback loop for Core's self-improvement, allowing Guardian to detect when Core's "brain" needs tuning.

**Security Considerations for Core-Guardian Interactions:**

- **Mutual Authentication:** Crucially, the gRPC channel between Core and Guardian for `ActionProposal` and `SystemConfigUpdate` uses mandatory mTLS, meaning both Core and Guardian cryptographically verify each other's identity before any data is exchanged. This prevents impersonation.
- **Fine-grained Authorization:** Guardian's gRPC server implements strict RBAC. Core's service account (authenticated via Kubernetes/Keycloak integrated with Foundry) is only authorized to call `policy_enforce` and receive `system_config_update` for its own component, not for other agents or critical Foundry services.
- **Least Privilege:** Core runs with the absolute minimum Kubernetes RBAC permissions necessary for its operation. It cannot directly modify Guardian's policies or models.
- **Policy Audit Logs:** Every `ActionProposal` and Guardian's verdict is recorded in immutable, tamper-evident audit logs (persisted in Foundry's data lake), providing a complete history for compliance and forensic analysis.

- **Rate Limiting:** Guardian's `/policy_enforce` endpoint implements rate limiting to prevent Core from overwhelming it with proposals, even if Core malfunctions.
- **Content Sanitization (for Proposals):** While `ActionProposal` data is primarily structured, any embedded text (like `context_summary` or `reason` fields) is sanitized by Guardian to prevent injection attacks or data corruption.

Script Kitty.Foundry

Script Kitty: The Genesis of Orchestrated Intelligence - Microscopic Dissection

This truly exhaustive, incredibly in-depth, expert-level technical analysis of Script Kitty's architecture is structured into **six distinct parts**:

1. **Script Kitty.Nexus**: The User's Portal & Group Chat Facilitator
2. **Script Kitty.Core**: The Central Governing AI / High-Level Orchestrator
3. **Script Kitty.Armory**: Resource Acquisition & Tooling AI
4. **Script Kitty.Skills**: Task Execution AI
5. **Script Kitty.Guardian**: Safety & Alignment Proxy / Evaluation & Training AI
6. **Script Kitty.Foundry**: MLOps & Infrastructure Fabric

Each part meticulously dissects every function, component, and tool, leaving absolutely no stone unturned, no parameter unscrutinized, to ensure flawless execution and the retention of all previously discussed abilities, capabilities, and core principles.

---

## Part 6: Script Kitty.Foundry - The MLOps & Infrastructure Fabric

**Core Purpose:** Script Kitty.Foundry is the bedrock, the unseen yet indispensable engine that propels the entire Script Kitty ecosystem. It is the sophisticated MLOps (Machine Learning Operations) and infrastructure orchestration layer, providing universal services for compute, data, security, deployment, and operational intelligence. Foundry ensures seamless transitions, continuous evolution, universal agent features, and a platform that learns and adapts with the collective knowledge of all Script Kitty agents, making the abstract notion of "AGI" a robust and practical reality. It underpins every single function, component, and interaction within Script Kitty, guaranteeing its unparalleled functionality and relentless improvement.

**Technical Stack Overview (Foundry):**

- **Orchestration:** Kubernetes (the foundational distributed operating system for Script Kitty).
- **Containerization:** Docker.
- **CI/CD & GitOps:** GitLab CI/CD/GitHub Actions (CI/CD Pipelines), Argo CD/Flux CD (GitOps).
- **Observability:** Prometheus (metrics), Grafana (dashboards, alerts), Loki/Elasticsearch (logging), Fluent Bit/Fluentd (log collection), OpenTelemetry (tracing), Jaeger/Zipkin (trace visualization).
- **Data Management:** MinIO/Ceph/AWS S3/GCP GCS (Object Storage), DVC (Data Version Control), Feast (Feature Store), Apache Airflow/Prefect/Dagster (Data Orchestration).

- **Security:** HashiCorp Vault (secrets), Calico/Cilium (network policy), Istio/Linkerd (service mesh), Falco (runtime security), Trivy/Clair (container scanning), Keycloak (identity management).
- **ML Pipelines:** Kubeflow Pipelines/Argo Workflows (workflow orchestration), MLflow (experiment tracking), Ray Tune/Optuna (hyperparameter optimization).
- **Virtualization (Hybrid):** KubeVirt (for specialized legacy workloads if needed).

---

**Minute Aspects & Technical Dissection:**

**1. Kubernetes-Native Compute & Orchestration Sub-system**

This is the very operating system of Script Kitty's distributed brain. It is responsible for the deployment, scaling, healing, and resource management of every microservice and specialized agent.

- **1.1. Kubernetes Cluster Management (Control Plane & Worker Nodes):**

  - **Purpose:** Provides the declarative, self-healing, and scalable foundation for all Script Kitty operations.
  - **Technical Dissection:**
    - **API Server:** The central control plane component, exposing the Kubernetes API (RESTful interface) for all interactions. Every deployment, scaling event, or configuration change for Script Kitty's agents goes through this API. Authentication and authorization (via RBAC - Role-Based Access Control) are meticulously enforced for all API calls.
    - **etcd:** The highly available key-value store, serving as Kubernetes's consistent and persistent backbone for cluster state and configuration. All desired states for Script Kitty's services are stored here.
    - **Scheduler:** Monitors newly created pods and assigns them to worker nodes based on resource requirements (CPU, memory, GPU), node constraints (taints and tolerations), and affinity/anti-affinity rules. Foundry's scheduler ensures optimal placement of Script Kitty's diverse agents (e.g., placing GPU-intensive LLMs on GPU nodes).
    - **Controller Manager:** Runs various controllers (e.g., Deployment Controller, ReplicaSet Controller) that continuously watch the state of the cluster and make changes to move the current state towards the desired state. For Script Kitty, this means ensuring the correct number of Nexus, Core, Armory, Skills, and Guardian pods are always running.
    - **Kubelet (on Worker Nodes):** The agent that runs on each worker node, ensuring containers are running in a Pod. It communicates with the API server, manages Pods, mounts volumes, and reports node status. Kubelet is critical for executing Script Kitty's agent workloads.

- **Kube-proxy (on Worker Nodes):** Maintains network rules on nodes, enabling network communication to your Pods from inside or outside the cluster. It ensures that Script Kitty's internal microservices can communicate seamlessly and reliably.
- **Horizontal Pod Autoscaler (HPA):** Dynamically scales the number of running pods for stateless services (e.g., Nexus Channel Adapters, LLM inference services, Core's message processors) based on observed CPU utilization, memory consumption, or custom metrics (e.g., inbound message queue depth, active user sessions). This ensures Script Kitty can handle fluctuating loads without manual intervention.
- **Vertical Pod Autoscaler (VPA) (Experimental/Advisory):** While primarily used for recommendations, VPA can adjust resource requests/limits for pods, optimizing resource allocation based on historical usage. This helps right-size Script Kitty's specialized agents over time.
- **Cluster Autoscaler:** Automatically adjusts the number of worker nodes in the Kubernetes cluster based on pending pods and resource utilization. This is critical for scaling Script Kitty's entire compute fabric up and down dynamically with demand, optimizing cloud costs.
- **Node Affinity/Anti-Affinity & Taints/Tolerations:** Used extensively to ensure specific agent types (e.g., GPU-intensive LLMs for Core/Skills, high-I/O data processing for Armory) are scheduled on appropriate nodes, or to prevent co-location of sensitive services on the same node for security/performance reasons.

- **1.2. Hybrid Compute Capabilities (KubeVirt Integration - Conditional):**

  - **Purpose:** While primarily container-native, Foundry can optionally integrate KubeVirt to run virtual machines alongside containers within the same Kubernetes cluster. This provides flexibility for incorporating legacy tools or specific software that absolutely requires a VM environment (e.g., highly specialized proprietary quantum computing simulators, older enterprise software).
  - **Technical Dissection:**
    - **KubeVirt CRDs:** Custom Resource Definitions (CRDs) like `VirtualMachine` are installed, allowing VMs to be managed like any other Kubernetes resource.
    - **VM Lifecycle Management:** KubeVirt controllers manage the lifecycle of VMs (start, stop, pause, migrate) using standard Kubernetes commands.
    - **Network Integration:** VMs share the same Kubernetes network fabric, allowing seamless communication with containerized services.
    - **Persistent Volumes:** VMs can utilize Kubernetes Persistent Volumes, integrating with Foundry's storage management.
    - **Use Cases for Script Kitty:** Potentially for specialized "Skills" agents that rely on non-containerizable software, or for a highly isolated execution environment for extremely sensitive "Armory" tool wrappers.

- **1.3. Custom Resource Definitions (CRDs) for Script Kitty Agents:**

  - **Purpose:** Foundry defines custom Kubernetes resources that represent Script Kitty's specialized agents (e.g., `ScriptKittyAgent`, `SkillModule`, `ToolManifest`). This allows Kubernetes to understand and manage these high-level logical constructs directly.
  - **Technical Dissection:**
    - **YAML Definitions:** CRDs are defined in YAML files, specifying their schema (`spec` and `status`) and versioning.
    - **Custom Controllers:** Foundry implements custom Kubernetes controllers (often using operator frameworks like Kubebuilder or Operator SDK) that watch for changes to these CRDs. When a `ScriptKittyAgent` CR is created or updated, the custom controller translates this into standard Kubernetes Deployments, Services, and other resources.
    - **Declarative Agent Management:** This enables Core or even human operators to declaratively define desired agent configurations (e.g., "deploy an instance of CodeLlama-7B skill module with 2 GPUs") via CRs, and Kubernetes handles the underlying orchestration.
    - **Status Reporting:** The custom controller updates the `status` field of the CRD, reflecting the real-time operational status of the agent (e.g., `Running`, `Degraded`, `Updating`).

## 2. Container Registry & Image Management Sub-system

The secure repository for all of Script Kitty's executable components, ensuring versioning, integrity, and controlled distribution.

- **2.1. Secure Container Registry:**
  - **Purpose:** Stores all Docker images for Script Kitty components (Nexus, Core, Guardian, Foundry services) and agent models (Armory's tool wrappers, Skills' specialized ML models).
  - **Technical Dissection:**
    - **Choice of Registry:**
      - **Self-hosted:** Harbor or Quay.io for complete control, enhanced security features (Vulnerability Scanning, Content Trust, Replication), and integration with internal networks.
      - **Cloud-Managed:** AWS ECR, GCP GCR, Azure Container Registry for ease of management, high availability, and integration with cloud-native security services.
    - **Image Tagging & Versioning:** Strict adherence to semantic versioning (e.g., `v1.2.3`, `v1.2.3-commitsha`, `latest-stable`) for all images, facilitating rollbacks and ensuring traceability.

- **Access Control:** Robust RBAC (Role-Based Access Control) integrated with Script Kitty's central identity management (Keycloak) to strictly control who can push, pull, or delete images.
- **Image Signing (Content Trust):** Utilizes Notary (or equivalent) for cryptographically signing container images, verifying their authenticity and integrity before deployment to prevent tampering or supply chain attacks.
- **Vulnerability Scanning:** Integrated scanners (e.g., Trivy, Clair, Harbor's built-in scanner) automatically scan all pushed images for known CVEs (Common Vulnerabilities and Exposures), flagging insecure dependencies and preventing deployment of vulnerable components.
- **Image Immutability:** Once an image is tagged and pushed, it is immutable, ensuring that the deployed artifact is precisely what was built and tested.

## 3. CI/CD & GitOps Automation Sub-system

The core of Script Kitty's continuous evolution, automating the entire software delivery lifecycle with a strong emphasis on declarative configuration and version control.

- **3.1. Continuous Integration (CI) Pipelines:**

  - **Purpose:** Automates the build and test phases for every code change made across Script Kitty's repositories.
  - **Technical Dissection:**
    - **Triggering:** Pipelines are automatically triggered by Git events (e.g., `push` to `main` branch, `pull_request` creation).
    - **Build Jobs:** Compiles source code (e.g., Python, Go), resolves dependencies, and builds Docker images for microservices and agent models.
    - **Unit Tests:** Executes comprehensive unit tests for individual functions and classes.
    - **Integration Tests:** Validates interactions between closely related components (e.g., Nexus NLU with Session Manager).
    - **Linting & Static Analysis:** Enforces coding standards and detects potential bugs or security vulnerabilities (e.g., `Pylint`, `ESLint`, `Bandit`, `SonarQube` community edition).
    - **Security Scans (Pre-Push):** Runs basic vulnerability scans on dependencies and code before image creation.
    - **Artifact Generation:** If all tests pass, it generates deployable artifacts (e.g., Docker images, Helm charts) and pushes them to the secure container registry with appropriate version tags.
    - **Tools:** GitLab CI/CD, GitHub Actions, Jenkins, Tekton Pipelines are all viable orchestrators for these complex, multi-stage pipelines.

- **3.2. Continuous Delivery/Deployment (CD) & GitOps:**

  - **Purpose:** Automates the deployment and update processes for all Script Kitty services and agents, ensuring that the live system always reflects the desired state defined in Git.
  - **Technical Dissection:**
    - **Declarative Infrastructure as Code (IaC):** All Kubernetes manifests (Deployments, Services, ConfigMaps, CRDs, Helm Charts) are stored in a Git repository, serving as the single source of truth for Script Kitty's infrastructure and application configurations.
    - **GitOps Agents:** Argo CD or Flux CD are deployed within the Kubernetes cluster. These agents continuously monitor the Git repositories for changes.
    - **Automated Synchronization:** When a change is detected in Git (e.g., a new Docker image tag for Nexus), the GitOps agent automatically pulls the new configuration and applies it to the Kubernetes cluster, reconciling the actual state with the desired state.
    - **Deployment Strategies:**
      - **Rolling Updates (Default):** Kubernetes natively supports rolling updates, gradually replacing old pods with new ones.
      - **Blue/Green Deployments:** For critical services, Foundry orchestrates Blue/Green deployments (via GitOps or custom controllers). A new "Green" environment is deployed alongside the existing "Blue" one. Once verified, traffic is instantly shifted to Green, minimizing downtime. Rapid rollback to Blue is possible.
      - **Canary Deployments:** A small percentage of user traffic is routed to new "Canary" version of a service. Foundry's monitoring (via Guardian's integration) meticulously tracks metrics (latency, error rates, user feedback) for the Canary. If performance is stable, traffic is gradually shifted; otherwise, the Canary is rolled back. This is crucial for safely deploying new LLM models or complex agent logic.
    - **Automated Rollbacks:** If a deployment fails (e.g., pods crash, health checks fail, or Guardian detects anomalies post-deployment), the GitOps agent automatically rolls back to the last known good configuration in Git.
    - **Immutable Deployments:** Once a container image is built and tested, it's deployed as is. Configuration changes are applied via ConfigMaps or Secrets, which trigger rolling updates.

## 4. Observability Stack (Centralized Logging, Metrics, Tracing) Sub-system

The eyes and ears of Script Kitty, providing real-time insights into its operational health, performance, and behavior. Essential for debugging, monitoring, and enabling continuous improvement.

- **4.1. Metrics Collection & Monitoring (Prometheus & Grafana):**

  - **Purpose:** Gathers high-frequency numerical data about the system's state and performance.
  - **Technical Dissection:**
    - **Prometheus:** A powerful time-series database and monitoring system. Each Script Kitty microservice (Nexus, Core, Armory, Skills, Guardian, Foundry components) exposes a `/metrics` endpoint in the Prometheus text format, containing granular metrics (e.g., `http_requests_total`, `api_call_latency_seconds`, `llm_inference_duration_milliseconds`, `gpu_utilization_percent`, `agent_task_success_rate`).
    - **Prometheus Operator:** Deployed in Kubernetes, it automatically discovers and scrapes metrics from all Script Kitty services based on Kubernetes labels/annotations.
    - **Grafana:** A visualization and analytics platform. Pre-built and custom dashboards provide real-time views of Script Kitty's health (e.g., overall system latency, individual agent error rates, resource saturation, active user sessions, task completion rates).
    - **PromQL:** Complex queries are written using PromQL (Prometheus Query Language) to analyze trends, aggregate data, and identify specific performance bottlenecks within Script Kitty.
    - **Alerting:** Prometheus Alertmanager integrates with Grafana to trigger critical alerts (e.g., "Core planning failures exceeding threshold," "Skills sandboxed execution timeout," "Guardian policy violations detected") via Slack, PagerDuty, or email.
- **4.2. Centralized Logging (Fluent Bit/Fluentd & Loki/Elasticsearch):**

  - **Purpose:** Aggregates and stores all log messages generated by Script Kitty's components, enabling debugging, auditing, and post-mortem analysis.
  - **Technical Dissection:**
    - **Log Collection Agents:** Fluent Bit (lightweight, high-performance) or Fluentd are deployed as DaemonSets on each Kubernetes node. They collect container logs (from stdout/stderr) and node-level logs.
    - **Log Aggregation:**
      - **Loki:** Preferred for its Prometheus-inspired query language (LogQL) and efficient storage of semi-structured logs. Ideal for troubleshooting, enabling fast searches for specific `session_id`, `agent_name`, or `error_message`.
      - **Elasticsearch/OpenSearch:** For more advanced full-text search, complex aggregations, and long-term archival of highly structured logs.

- **Log Enrichment:** Logs are automatically enriched with Kubernetes metadata (pod name, namespace, container ID) and Script Kitty-specific identifiers (session ID, task ID) for easier filtering and correlation.
- **Log Retention Policies:** Configurable policies ensure efficient storage management, balancing compliance requirements with cost (e.g., hot storage for recent logs, cold storage for older logs).

- **4.3. Distributed Tracing (OpenTelemetry & Jaeger/Zipkin):**

  - **Purpose:** Provides end-to-end visibility into the flow of a single request or task execution across multiple microservices within Script Kitty. Crucial for understanding latency, identifying bottlenecks, and debugging complex distributed interactions.
  - **Technical Dissection:**
    - **OpenTelemetry SDKs:** Integrated into the code of every Script Kitty microservice (Nexus, Core, Armory, Skills, Guardian, Foundry's own services). These SDKs automatically instrument incoming/outgoing requests, database calls, and function executions.
    - **Span Generation:** Each operation within a request (e.g., "Nexus receives query," "Core plans task," "Armory fetches tool," "Skills executes code") generates a "span." Spans contain metadata (start/end time, duration, attributes like `http.method`, `db.statement`, `agent.name`, `task.id`).
    - **Trace Context Propagation:** Unique `trace_id` and `span_id` are propagated across service boundaries (e.g., HTTP headers, gRPC metadata), linking all related spans into a single "trace."
    - **Trace Exporters:** Spans are exported to a backend collector (e.g., OpenTelemetry Collector).
    - **Jaeger/Zipkin:** Backend systems for storing, indexing, and visualizing traces. Developers and SREs can use Jaeger UI to visualize a flame graph of a single user request, showing which services were involved, how long each step took, and where errors occurred.
    - **Correlation with Logs/Metrics:** Trace IDs are often included in log messages, allowing direct navigation from a trace view to relevant log entries.

## 5. Data Management (Data Lake, Data Versioning, Feature Store) Sub-system

The robust infrastructure for all data assets within Script Kitty, from raw input to trained models, ensuring integrity, accessibility, and reproducibility.

- **5.1. Scalable Object Storage (Data Lake Foundation):**

  - **Purpose:** Provides durable, cost-effective, and massively scalable storage for all unstructured and semi-structured data generated or consumed by Script Kitty.
  - **Technical Dissection:**

- 
    - 
        - **Choice of Storage:**
            - **Self-hosted:** MinIO (S3-compatible object storage, can run on Kubernetes) or Ceph (distributed object, block, and file storage for large-scale deployments).
            - **Cloud-Managed:** AWS S3, GCP GCS, Azure Blob Storage (for unparalleled scalability, durability, and integration with other cloud services).
        - **Data Types:** Stores raw user inputs, web scraped data, extracted research documents, model training datasets, pre-processed features, model checkpoints, inference logs, audit logs, and Guardian's feedback data.
        - **Bucket/Prefix Organization:** Data is logically organized into buckets (or prefixes within a single bucket) based on data source, type, and stage of processing (e.g., `raw-user-inputs`, `armory-research-docs`, `skills-model-checkpoints`, `guardian-feedback-loop`).
        - **Lifecycle Policies:** Configured for automatic tiering (e.g., moving older data from hot to cold storage) and deletion, optimizing storage costs.
        - **Access Control:** Granular IAM policies control read/write access to specific buckets or prefixes for different Script Kitty agents.
- **5.2. Data Version Control (DVC - Data Version Control):**

    - **Purpose:** Versions datasets and ML models like source code, ensuring reproducibility of experiments and deployments.
    - **Technical Dissection:**
        - **Git Integration:** DVC works alongside Git. Git tracks small `.dvc` files that point to the actual data files, which are stored in object storage.
        - **Data Checksums:** DVC generates checksums for data files, allowing it to detect changes.
        - **Pipeline Definition:** DVC can define data processing and ML training pipelines (`dvc.yaml`), linking code, data, and models. This is critical for Script Kitty's continuous learning. For example, Guardian triggering a retraining of a Skills agent model will use DVC to fetch the exact version of the training data and code.
        - **Reproducibility:** Any team member can `dvc checkout` a specific version of the data and code from Git, ensuring they can reproduce exact training runs or model evaluations.
- **5.3. Feature Store (Feast):**

    - **Purpose:** Provides a centralized, standardized, and discoverable repository for machine learning features, ensuring consistency between training and inference environments.
    - **Technical Dissection:**

- **Offline Store:** Typically uses a data lake (S3/GCS) or a data warehouse (Snowflake/BigQuery) for historical feature data used in model training by Guardian.
- **Online Store:** Uses a low-latency key-value store (e.g., Redis, Cassandra) for real-time feature serving during model inference by Skills agents (e.g., fetching a user's recent activity score for a personalized response, or a tool's success rate for Armory's decision making).
- **Feature Definition:** Features are defined declaratively in `feature_store.yaml` files, including data types, sources, and transformation logic.
- **Point-in-Time Correctness:** Feast ensures that historical features used for training accurately reflect the state of the world at the time the events occurred, preventing data leakage.
- **Use Cases for Script Kitty:**
  - **Skills:** Provides pre-computed features for specialized AI models (e.g., historical user interaction patterns for a recommendation model).
  - **Guardian:** Stores performance metrics and feedback signals as features for models that predict optimal retraining schedules or detect anomalies.
  - **Core:** Might use features related to past task success rates for its planning engine.

- **5.4. Data Orchestration (Apache Airflow/Prefect/Dagster):**

  - **Purpose:** Manages and schedules complex data pipelines for ingestion, transformation, and movement of data within Script Kitty's data lake.
  - **Technical Dissection:**
    - **DAGs (Directed Acyclic Graphs):** Data pipelines are defined as DAGs in Python code, representing sequences of tasks with dependencies.
    - **Scheduling:** Tasks can be scheduled at fixed intervals (e.g., daily ingestion of new web data for Armory) or triggered by events (e.g., new model artifact pushed).
    - **Error Handling & Retries:** Robust mechanisms for handling task failures, with configurable retries, alerting, and manual intervention points.
    - **Monitoring & UI:** Provides a web UI for monitoring pipeline runs, viewing logs, and manually triggering tasks.
    - **Use Cases for Script Kitty:**
      - **Armory:** Orchestrates web scraping, data cleaning, and data loading into the data lake.
      - **Guardian:** Manages pipelines for collecting feedback data, generating evaluation metrics, and preparing datasets for model retraining.
      - **Foundry's Internal Data:** Processes logs and metrics for long-term analytics.

**6. Security & Identity Management Sub-system**

The paramount layer ensuring the integrity, confidentiality, and availability of Script Kitty, protecting it from external threats and internal misuse.

- **6.1. Secrets Management (HashiCorp Vault):**

  - **Purpose:** Securely stores, accesses, and manages sensitive credentials (API keys, database passwords, cloud tokens) required by Script Kitty's services and agents.
  - **Technical Dissection:**
    - **Centralized Storage:** All secrets are stored in an encrypted, audited, and versioned central repository.
    - **Dynamic Secrets:** Vault can generate on-demand, short-lived credentials for databases, cloud providers, and other services, minimizing the risk of long-lived, static secrets. For example, a Skills agent needing to access a database will request a temporary credential from Vault.
    - **Leasing & Revocation:** Secrets have associated leases and can be automatically revoked after use or expiry.
    - **Audit Logging:** All access to secrets is meticulously logged for auditing and compliance.
    - **Authentication & Authorization:** Integrates with Kubernetes Service Accounts and other identity providers (e.g., Keycloak) to authenticate clients requesting secrets and authorize access based on policies.
- **6.2. Network Policy (Calico/Cilium):**

  - **Purpose:** Enforces network segmentation and traffic filtering within the Kubernetes cluster, creating a "zero-trust" environment where agents can only communicate with explicitly allowed services.
  - **Technical Dissection:**
    - **Declarative Rules:** Network policies are defined as Kubernetes YAML manifests, specifying allowed ingress (inbound) and egress (outbound) traffic based on labels, namespaces, IP blocks, and ports.
    - **Micro-segmentation:** Enforces isolation between Script Kitty's services (e.g., Nexus can talk to Core, but Skills cannot directly talk to Core's sensitive internal endpoints unless explicitly allowed).
    - **Prevention of Lateral Movement:** Restricts an attacker from moving laterally across the cluster if one service is compromised. For example, a compromised Skills agent cannot initiate connections to the Guardian's policy enforcement engine or Foundry's secrets management.
    - **DNS Policies:** Can restrict DNS lookups to authorized internal services only.
- **6.3. Service Mesh (Istio/Linkerd - Security Features):**

- ○ **Purpose:** Adds a layer of sophisticated traffic management, observability, and *enhanced security* capabilities to Script Kitty's inter-service communication without modifying application code.
- ○ **Technical Dissection:**
    - ■ **mTLS (Mutual TLS):** Automatically encrypts and authenticates all traffic between Script Kitty's microservices, ensuring that only trusted services can communicate. This is a crucial defense against eavesdropping and impersonation.
    - ■ **Fine-grained Access Control:** Policy enforcement at the network layer (e.g., "only Nexus can call Core's `process_query` endpoint, and only if authenticated"). This goes beyond simple network policies by understanding service identity.
    - ■ **Authorization Policies:** Define granular `AuthorizationPolicy` rules (e.g., `allow` or `deny` specific HTTP methods, paths, or source services).
    - ■ **Traffic Mirroring/Shadowing:** Allows for mirroring production traffic to a new version of a service for testing purposes without impacting live users, valuable for safely testing new agent behaviors.
- ● **6.4. Runtime Security (Falco):**

    - ○ **Purpose:** Detects anomalous and suspicious behavior at the kernel level within Script Kitty's running containers, providing real-time threat detection.
    - ○ **Technical Dissection:**
        - ■ **System Call Monitoring:** Falco monitors system calls (`syscalls`) directly from the Linux kernel, providing deep visibility into container activity.
        - ■ **Rule Engine:** Uses a declarative rule language to define suspicious behaviors (e.g., a process attempting to write to `/etc/passwd`, a container spawning an unexpected shell, a Skills agent trying to access unauthorized files, network connections to unknown IPs).
        - ■ **Real-time Alerts:** Generates alerts (e.g., to Guardian, to a SIEM system, or to a Slack channel) when a rule is violated, enabling immediate response to potential compromises.
        - ■ **Use Cases for Script Kitty:** Detects sandboxed escape attempts by Skills agents, unauthorized data exfiltration by Armory, or suspicious administrative activity within the Core or Foundry components.
- ● **6.5. Container Scanning (Trivy/Clair):**

    - ○ **Purpose:** Scans Docker images for known vulnerabilities and misconfigurations before deployment.
    - ○ **Technical Dissection:**
        - ■ **Layer-by-Layer Scanning:** Analyzes each layer of a Docker image for vulnerabilities in operating system packages, language-specific dependencies (Python, Go, Node.js), and application libraries.

- ■ **Vulnerability Databases:** Integrates with public CVE databases (NVD) and commercial vulnerability intelligence feeds.
- ■ **SBOM (Software Bill of Materials) Generation:** Can generate a list of all components and dependencies within an image, improving transparency and auditability.
- ■ **Integration with CI/CD:** Scans are integrated into CI pipelines. Images with high-severity vulnerabilities can automatically fail the build, preventing their deployment.
- ■ **Use Cases for Script Kitty:** Ensures all base images and application dependencies for Nexus, Core, Armory, Skills, and Guardian are regularly scanned and patched.

- ● **6.6. Identity & Access Management (Keycloak):**

  - ○ **Purpose:** Provides centralized authentication and authorization services for human users interacting with Script Kitty's administrative UIs (e.g., Guardian's oversight dashboard, Foundry's MLOps dashboards) and for internal service-to-service authentication.
  - ○ **Technical Dissection:**
    - ■ **OpenID Connect/OAuth2:** Implements standard protocols for secure authentication and authorization.
    - ■ **User/Role Management:** Manages users, groups, and roles, defining what actions users or internal services are permitted to perform.
    - ■ **Single Sign-On (SSO):** Enables human operators to log in once and access multiple Script Kitty administrative applications securely.
    - ■ **Service Account Management:** Provides a robust way to manage credentials and access policies for internal Script Kitty service accounts (e.g., Foundry's CI/CD pipelines accessing Kubernetes API, Guardian's feedback system updating the GKGS).
    - ■ **Integration with Kubernetes:** Service accounts are typically linked to Keycloak roles, allowing Kubernetes to automatically authorize internal pod communication based on their assigned identity.

---

**Foundry's Role in Script Kitty's Overall Functionality:**

- ● **"Seamless Transitions":** Achieved through robust CI/CD pipelines, GitOps-driven deployments (Blue/Green, Canary), and Kubernetes's rolling update capabilities, ensuring new models and features are deployed without disruption. Observability provides the confidence to make these transitions.
- ● **"Evolving":** Foundry is the enabler of continuous learning and improvement. It provides the data versioning, ML pipeline orchestration, and model serving infrastructure that Guardian uses to trigger and manage retraining cycles, and that Core uses to deploy new planning models.

- **"Universal Agent Features":** Every Script Kitty agent (Nexus, Core, Armory, Skills, Guardian) inherently benefits from Foundry's services: containerization, secure communication (mTLS via service mesh), centralized logging, metrics collection, secrets management, and data access. Foundry provides the standardized "platform" upon which all agents are built.
- **"Platform that evolves with all agents' knowledge":** Foundry provides the data infrastructure (Data Lake, DVC, Feature Store) for accumulating the collective experience and knowledge (e.g., Core's successful plans, Armory's extracted data, Skills' task outcomes, Guardian's feedback). This accumulated data then feeds back into retraining cycles orchestrated by Foundry's ML pipelines, leading to the evolution of models across the entire Script Kitty system. The MLOps practices ensure that this knowledge is consistently applied and measured.

# Foundry Interactions

## Part 7: Inter-Component Interactions - Foundry as the Foundational Nexus

Foundry is not merely a component; it is the **omnipresent substrate**, the **nervous system's backbone**, and the **lifeblood** of Script Kitty. It is the sophisticated MLOps and infrastructure fabric that *enables* all other components to exist, operate, scale, communicate, and evolve. Its interactions are not merely data exchanges but deep symbiotic relationships, providing the very operational environment and the mechanisms for continuous adaptation that allow Script Kitty to function as a cohesive AGI. We begin by meticulously dissecting the interactions between Foundry and every other component.

---

### 1. Foundry's Foundational Role & Interaction Philosophy

Foundry's core philosophy is to provide **universal, abstracted, and declarative services** that all other Script Kitty components consume. It doesn't typically initiate complex "requests" to other components in the same way Core might request Armory to research. Instead, Foundry provides:

- **Compute Orchestration:** Provisioning and managing the Kubernetes environment where *all* other components run.
- **Deployment & Lifecycle Management:** Automating the build, deployment, scaling, and updates of Nexus, Core, Armory, Skills, and Guardian microservices and models.
- **Observability Fabric:** Ingesting metrics, logs, and traces from every component, providing a unified operational view.
- **Data Infrastructure:** Offering scalable storage, versioning, and feature serving that components leverage.
- **Security Perimeter:** Enforcing identity, access control, network policies, and runtime security for the entire ecosystem.
- **MLOps Automation:** Providing the pipelines and frameworks for model training, evaluation, and deployment, which Guardian and Core utilize.

Essentially, Foundry is the **provider**, and the other components are **consumers** of its infrastructure and platform services.

---

### 2. Foundry ↔ Nexus: The Operational Foundation for User Interaction

The interaction between Foundry and Nexus is fundamental for Script Kitty's user-facing capabilities. Foundry provides Nexus with everything it needs to be alive, scalable, observable, and continuously updated.

- **2.1. Deployment & Lifecycle Management:**

- ○ **Foundry's Role:** Foundry's CI/CD & GitOps Automation sub-system (3.0) is directly responsible for building, packaging, and deploying all Nexus microservices.
- ○ **Microscopic Detail:**
  - ■ **Trigger:** A developer committing a change to the Nexus codebase (e.g., updating a channel adapter, refining the response synthesis logic, or updating the Nexus LLM model) triggers a GitLab CI/CD or GitHub Actions pipeline defined within Foundry's automation.
  - ■ **Build Process (Foundry's CI):** The pipeline orchestrates:
    - ■ **Code Compilation/Dependency Resolution:** `pip install -r requirements.txt`, `go mod download`.
    - ■ **Unit/Integration Tests:** Executes `pytest` for Python services, `go test` for Go services.
    - ■ **Container Image Build:** A `Dockerfile` (e.g., `FROM python:3.10-slim-bookworm`, `COPY . /app`, `RUN pip install -r requirements.txt`, `CMD ["python", "app.py"]`) defines the Nexus service. Foundry's CI tool (`docker build -t nexus-channel-adapter:v1.2.3 .`) builds the image.
    - ■ **Image Scanning:** `Trivy scan --severity CRITICAL,HIGH nexus-channel-adapter:v1.2.3` (Foundry's Container Scanning 6.5) runs automatically to detect vulnerabilities. If critical vulnerabilities are found, the build fails.
    - ■ **Image Push:** The built and scanned image is pushed to Foundry's Secure Container Registry (2.1), tagged with a unique version (e.g., `nexus-channel-adapter:v1.2.3-commitsha`).
  - ■ **Deployment Process (Foundry's GitOps):**
    - ■ **Git Commit:** The CI pipeline updates a Kubernetes manifest (e.g., a Helm chart's `values.yaml` or a plain Kubernetes `Deployment.yaml`) in a GitOps repository with the new image tag.
    - ■ **Argo CD/Flux CD Sync:** Foundry's GitOps Agent (3.2) continuously monitors this GitOps repository. Upon detecting the change, it pulls the new manifest.
    - ■ **Kubernetes API Interaction:** Argo CD/Flux CD interacts with Foundry's Kubernetes API Server (1.1) to apply the new `Deployment` specification for the Nexus service.
    - ■ **Rolling Update:** Kubernetes (managed by Foundry's Kubelet and Controller Manager 1.1) performs a rolling update, gradually replacing old Nexus pods with new ones. `minReadySeconds` and `maxUnavailable` parameters ensure smooth transitions, minimizing downtime.

- **Health Checks:** `livenessProbe` and `readinessProbe` are defined in the Nexus deployment. Foundry's Kubelet continuously checks these endpoints (e.g., HTTP GET to `/healthz`). If health checks fail, Kubernetes prevents traffic from being routed to the new pod and can trigger a rollback.
  - **Model Updates (Nexus LLM):** New versions of the Nexus LLM (fine-tuned or updated base model) are packaged into new Docker images for the `vLLM` or `Ollama` inference service. Foundry's CI/CD orchestrates the build, scan, and GitOps-driven deployment of this new LLM inference service, seamlessly swapping out the underlying model via the same rolling update mechanism.
- **2.2. Compute Resource Provisioning & Scaling:**
  - **Foundry's Role:** Provides Nexus with the necessary CPU, memory, and potentially GPU resources via Kubernetes and manages its dynamic scaling.
  - **Microscopic Detail:**
    - **Resource Requests/Limits:** Nexus deployments specify `resources.requests.cpu`, `resources.requests.memory`, and `resources.limits` in their Kubernetes manifests. Foundry's Scheduler (1.1) uses these to place pods. For Nexus LLM, `resources.requests.nvidia.com/gpu` (if GPU-accelerated) would be defined.
    - **Horizontal Pod Autoscaling (HPA):** Foundry's HPA (1.1) monitors Nexus services (e.g., `nexus-channel-adapter`, `nexus-llm-inference`).
      - **Metrics:** It scrapes metrics from Prometheus (4.1) like `http_requests_per_second_avg`, `cpu_utilization_percentage`, or `active_sessions_count_per_pod` (custom metric exposed by Nexus and scraped by Prometheus).
      - **Scaling Logic:** `apiVersion: autoscaling/v2beta2`, `kind: HorizontalPodAutoscaler`, `spec.minReplicas`, `spec.maxReplicas`, `spec.metrics.resource.target.averageUtilization` or `spec.metrics.external.target.value`. If, for example, `cpu_utilization_percentage` exceeds 70%, HPA will tell Kubernetes to add more Nexus Channel Adapter pods.
    - **Cluster Autoscaler:** If the cluster runs out of nodes to schedule new Nexus pods, Foundry's Cluster Autoscaler (1.1) automatically provisions new worker nodes in the underlying cloud provider, ensuring Nexus's scalability.
- **2.3. Observability & Monitoring:**

- ○ **Foundry's Role:** Collects, aggregates, and visualizes all operational data from Nexus, enabling real-time insights and proactive issue detection.
- ○ **Microscopic Detail:**
    - ■ **Metrics (Nexus to Foundry):**
        - ■ Nexus services (e.g., `nexus-channel-adapter`, `nexus-nlu`, `nexus-response-synthesis`) expose Prometheus `/metrics` endpoints.
        - ■ Foundry's Prometheus (4.1) (managed by Prometheus Operator) is configured to scrape these endpoints based on Kubernetes service labels/annotations (`scrape_configs: - job_name: 'nexus-adapters' kubernetes_sd_configs: - role: pod relabel_configs: - source_labels: [__meta_kubernetes_pod_label_app] action: keep regex: nexus-channel-adapter`).
        - ■ Metrics captured: `nexus_messages_received_total`, `nexus_nlu_latency_seconds_bucket`, `nexus_llm_inference_duration_seconds_sum`, `nexus_active_sessions_gauge`, `nexus_errors_total`.
        - ■ **Grafana Dashboards:** Foundry's Grafana (4.1) provides specialized dashboards for Nexus, visualizing these metrics, showing trends, and highlighting anomalies (e.g., sudden spikes in NLU latency, increase in error rates on a specific channel adapter).
    - ■ **Logging (Nexus to Foundry):**
        - ■ Nexus services log to `stdout`/`stderr` using structured logging libraries (e.g., Python's `logging` module configured with JSON formatters).
        - ■ Foundry's Fluent Bit/Fluentd (4.2) agents, deployed as DaemonSets on every Kubernetes node, capture these container logs.
        - ■ **Log Aggregation:** Fluent Bit forwards logs to Foundry's Loki or Elasticsearch (4.2). Each log entry is enriched with Kubernetes metadata (pod name, namespace, container ID) and Nexus-specific fields (e.g., `session_id`, `message_id`, `intent`).
        - ■ **Log Querying:** Operators can query logs in Loki (using LogQL: `{container="nexus-channel-adapter"} |~ "ERROR" | json | session_id="abc"`) or Elasticsearch/Kibana to diagnose issues in Nexus.
    - ■ **Distributed Tracing (Nexus to Foundry):**
        - ■ OpenTelemetry SDKs (4.3) are embedded in Nexus's code, instrumenting key operations (e.g., `nexus_message_received`,

`nexus_nlu_processing`, `nexus_core_request`,
`nexus_response_generation`).

- **Trace Context Propagation:** Nexus ensures that `trace_id` and `span_id` are propagated through HTTP headers (for REST calls), gRPC metadata (for calls to Core), and internal message queues.
- Nexus traces are exported to Foundry's OpenTelemetry Collector, which then forwards them to Jaeger/Zipkin (4.3).
- **Trace Visualization:** Developers can use Jaeger UI to visualize a trace for a specific `session_id`, showing the exact execution path of a user's query through Nexus, including all internal component calls and their latencies, identifying bottlenecks (e.g., "LLM inference took 500ms").
- **Alerting:** Foundry's Prometheus Alertmanager (4.1) is configured with alert rules based on Nexus metrics (e.g., `ALERT NexusHighErrorRate IF rate(nexus_errors_total[5m]) > 10`). These alerts are dispatched to relevant teams via Slack, PagerDuty, or email, ensuring immediate attention to Nexus-related operational issues.

- **2.4. Security & Access Control:**
  - **Foundry's Role:** Provides Nexus with a secure operating environment, manages its network access, and secures its sensitive data.
  - **Microscopic Detail:**
    - **Container Runtime Security:** Nexus containers run on Foundry's Kubernetes nodes, which are monitored by Falco (6.4). If a Nexus container attempts a suspicious syscall (e.g., trying to write to `/dev/mem` or executing an unauthorized binary), Falco triggers an alert, notifying Guardian and security teams.
    - **Network Policy (Foundry's Calico/Cilium 6.2):** Foundry defines strict Kubernetes Network Policies for Nexus.
      - **Egress:** Nexus's channel adapters are permitted to make outbound connections only to trusted external messaging platform APIs (e.g., `slack.com:443`, `api.telegram.org:443`). Egress to other internal Script Kitty services is restricted to `Core` (e.g., `core-grpc-service:50051`) and Foundry's `Prometheus`, `Loki`, `Jaeger` endpoints. All other outbound traffic is denied by default.
      - **Ingress:** Nexus's channel adapters are allowed ingress only from their respective external webhooks or internal load balancers. Internal Nexus microservices (e.g., `nexus-nlu`) are allowed ingress only from `nexus-channel-adapter` or other authorized internal Nexus components.
    - **Secrets Management (Foundry's HashiCorp Vault 6.1):**

- Nexus doesn't hardcode API keys for messaging platforms. Instead, it uses Vault. Nexus pods are configured with a Kubernetes Service Account, which Vault authenticates.
- A Nexus pod's initialization container or sidecar (e.g., Vault Agent Injector) requests short-lived credentials (e.g., `SLACK_BOT_TOKEN`) from Vault using `vault kv get secret/nexus/slack_bot_token`. These tokens are injected as environment variables or mounted as temporary files.
- This prevents sensitive data from residing in container images or Git repositories.
- **Image Security:** Nexus images are regularly scanned by Foundry's Trivy/Clair (6.5) during CI/CD to prevent known vulnerabilities from being deployed.

---

### 3. Foundry ↔ Core: Enabling the Strategic Brain's Operations

Core, the strategic brain of Script Kitty, relies heavily on Foundry for its operational environment, planning execution, and integration with the broader system.

- **3.1. Core Deployment & Scaling:**
  - **Foundry's Role:** Similar to Nexus, Foundry manages Core's deployment, scaling, and updates via GitOps and Kubernetes HPA/VPA.
  - **Microscopic Detail:**
    - Core's primary component, the `PlanningEngine` (via LLM for high-level planning) is resource-intensive. Foundry ensures it runs on nodes with sufficient CPU/memory and potentially dedicated GPUs for a larger LLM (e.g., Llama 3-70B served by vLLM).
    - `resource.requests` and `resource.limits` are carefully set. `Horizontal Pod Autoscaler` scales Core's stateless components (e.g., message ingestion, agent orchestration listeners) based on internal queue depth or CPU utilization.
    - Stateful components (e.g., parts of the Goal Decomposition engine that might need local disk for temporary planning state, although most should be stateless) might utilize Kubernetes `StatefulSets` and `PersistentVolumes` managed by Foundry's storage capabilities.
- **3.2. Global Knowledge Graph (GKGS) & Context Store Infrastructure:**
  - **Foundry's Role:** Provides the underlying database and storage infrastructure for Core's GKGS.
  - **Microscopic Detail:**
    - **Database Provisioning:** Foundry uses Infrastructure as Code (IaC) templates (Terraform/Ansible 3.4 in Armory, orchestrated by Core, but Foundry manages the *runtime environment* for these IaC tools) to

provision the database clusters (Neo4j or ArangoDB) for the GKGS. This involves Kubernetes Operators for these databases (e.g., `neo4j-operator`, `arangodb-operator`) running within Foundry's Kubernetes.

- **Persistent Storage:** The GKGS database requires persistent, highly available storage. Foundry provides `PersistentVolumeClaims` (PVCs) which provision `PersistentVolumes` (PVs) from its underlying storage solutions (e.g., network attached storage, distributed block storage like Ceph, or cloud provider persistent disks). This ensures the GKGS state is never lost even if Core pods are restarted.
- **Vector Database Integration:** For the semantic search capabilities of the GKGS, Foundry provisions and manages the vector database (Weaviate/Chroma). This often involves deploying its dedicated containerized services within Kubernetes and providing persistent storage for its vector indexes.

- **3.3. Internal Communication Hub & Message Broker Infrastructure:**
  - **Foundry's Role:** Provisions and manages the high-performance message broker (Kafka/Pulsar) and gRPC infrastructure for Core's internal communication.
  - **Microscopic Detail:**
    - **Kafka/Pulsar Cluster:** Foundry deploys and manages the Kafka or Pulsar cluster within Kubernetes, often using dedicated Operators (e.g., Strimzi for Kafka). This involves deploying Zookeeper (for Kafka), Kafka brokers, and ensuring persistent storage for message logs.
    - **gRPC Services:** Foundry's Kubernetes-native service discovery automatically registers Core's gRPC service endpoints. When Nexus or other agents (via Core) need to communicate with Core via gRPC, they use Kubernetes `Service` names (e.g., `core-grpc-service.scriptkitty-core-namespace.svc.cluster.local:50051`), which Foundry's Kube-proxy (1.1) resolves to the correct Core pod.
    - **Network Policies:** Foundry configures strict network policies (6.2) to allow only authorized internal services (Nexus, Armory, Skills, Guardian) to connect to Core's gRPC and Kafka endpoints.

- **3.4. Feedback & Learning Integration Infrastructure:**
  - **Foundry's Role:** Provides the streaming data platforms and orchestration tools that enable Core's feedback loop and learning cycles.
  - **Microscopic Detail:**
    - **Streaming Platforms:** Foundry deploys and manages Apache Flink or Apache Spark Streaming clusters within Kubernetes. These clusters consume feedback data from Kafka (populated by Guardian) for real-time processing and analysis by Core's learning components.
    - **MLflow/DVC Integration:** Foundry provides the persistent storage (object storage) for MLflow artifact tracking and DVC's data store, which

Core's planning model retraining components (orchestrated by Guardian, but running on Foundry) will use.

- ■ **RL Environment:** If Core uses Reinforcement Learning (RL) for refining planning policies, Foundry would provision the necessary distributed compute environment (e.g., a Ray cluster on Kubernetes) for running RLlib experiments.

---

**4. Foundry ↔ Armory: The Foundation for Resource & Tooling Provisioning**

Armory, Script Kitty's resource acquisition and tooling AI, relies on Foundry for secure tool execution environments, dynamic resource provisioning, and managing its internal tool registry.

- ● **4.1. Sandboxed Tool Execution Environments:**
  - ○ **Foundry's Role:** Provides the secure, isolated, and ephemeral environments where Armory can execute external tools (CLI, Python scripts, API wrappers).
  - ○ **Microscopic Detail:**
    - ■ **Ephemeral Pods/Jobs:** When Armory needs to run a web scraper or a generated API wrapper, it requests Core, which then dispatches the task to a specific Skill instance (potentially the `GenericToolExecutionSkill`). Foundry's Kubernetes-Native Compute (1.0) creates a new, short-lived Kubernetes `Job` or `Pod` for this specific execution.
    - ■ **Minimal Privileges:** These pods run with extremely restricted Service Accounts and Security Contexts (`runAsNonRoot: true`, `allowPrivilegeEscalation: false`).
    - ■ **Resource Limits:** Strict `resources.limits` are applied to prevent resource exhaustion from runaway scripts.
    - ■ **Network Policy (Egress):** The `NetworkPolicy` for these ephemeral pods (managed by Foundry's Calico/Cilium 6.2) is highly restrictive. They are typically allowed outbound connections *only* to the specific external URLs (e.g., target website for scraping, known API endpoint) or internal services (e.g., object storage for saving scraped data, Skill's internal communication for reporting results). All other outbound connections are denied.
    - ■ **Volume Mounts:** Only specific, read-only `PersistentVolumes` or ephemeral `emptyDir` volumes are mounted for temporary storage, preventing access to the host filesystem.
    - ■ **Kernel-level Security (Falco):** Foundry's Falco (6.4) monitors these sandboxed pods for any anomalous system calls that might indicate an escape attempt or malicious activity (e.g., opening a raw socket, trying to write to `/proc`).
- ● **4.2. Dynamic Tool Registry Infrastructure:**

- ○ **Foundry's Role:** Provides the database and search infrastructure for Armory's dynamic tool registry.
- ○ **Microscopic Detail:**
    - ■ **PostgreSQL/Elasticsearch Deployment:** Foundry deploys and manages the PostgreSQL database (for structured tool metadata) and Elasticsearch/OpenSearch cluster (for semantic search over tool capabilities) within Kubernetes. This involves deploying their respective Operators and managing their `PersistentVolumes`.
    - ■ **Scalability:** These databases are deployed for high availability (e.g., PostgreSQL replication) and horizontal scalability for Elasticsearch, managed by Foundry.
    - ■ **Data Integration:** Foundry's data orchestration tools (Airflow/Prefect 5.4) can be used to manage pipelines that ingest new tool manifests (e.g., from public API documentation or internal codebases) into this registry.
- ● **4.3. Resource Provisioning via IaC Templates Execution:**
    - ○ **Foundry's Role:** Provides the runtime environment for Armory's IaC execution, ensuring secure and audited provisioning of external resources.
    - ○ **Microscopic Detail:**
        - ■ **Ephemeral IaC Pods:** When Armory (via Core's planning) decides to provision a cloud VM or a specialized GPU instance, it doesn't directly run Terraform or Ansible. Instead, Core dispatches a task to a `ProvisioningSkill` (a type of Skills agent), which then requests Foundry to spin up an ephemeral pod containing the necessary IaC tools (e.g., `terraform` binary, `ansible-playbook`).
        - ■ **Cloud Provider SDKs/CLIs:** These ephemeral pods include the relevant cloud provider SDKs (e.g., `boto3` for AWS, `gcloud` CLI for GCP) and are configured with specific, temporary, and granular IAM roles/service accounts obtained from Foundry's Vault (6.1). This principle of least privilege ensures that the provisioning pod can *only* perform the exact requested action (e.g., create a specific VM type in a specific VPC, not delete entire accounts).
        - ■ **Auditing:** All actions performed by these IaC pods (which are executing on Foundry's compute) are meticulously logged and traced by Foundry's Observability Stack (4.0), providing a comprehensive audit trail to Guardian (via Kafka topic `foundry_resource_provisioning_logs`).
        - ■ **Quantum Computing Access:** For quantum computing, Foundry provides secure access to quantum hardware SDKs (Qiskit, Cirq, PennyLane) within sandboxed environments. Armory receives task descriptions (e.g., "submit quantum circuit X to IBM Quantum"), and Foundry's `QuantumExecutionSkill` executes the job, managing API keys via Vault and adhering to strict access policies.

**5. Foundry ↔ Skills: The Execution Engine's Platform**

Skills agents, Script Kitty's specialized task executors, are entirely dependent on Foundry for their existence, secure execution, data access, and model serving.

- **5.1. Specialized AI Models & Inference Serving:**
  - **Foundry's Role:** Provides the scalable and performant infrastructure for deploying and serving all Skills' specialized ML models.
  - **Microscopic Detail:**
    - **Model Serving Frameworks:** Foundry deploys and manages model serving frameworks like TorchServe, TensorFlow Serving, KServe (KNative Serving for ML), or Ray Serve within Kubernetes. Each Skills model (e.g., a computer vision model, a time-series predictor, a code generation model) is deployed as a separate microservice.
    - **GPU Allocation:** Foundry's Kubernetes Scheduler (1.1) precisely allocates GPU resources to these Skills model servers based on their `resources.requests` (e.g., `nvidia.com/gpu: 1`).
    - **Model Hot-Swapping/Versioning:** Foundry's CI/CD and GitOps (3.0) enable seamless updates to Skills models. New model versions (e.g., fine-tuned by Guardian) are deployed as new Docker images for the serving framework, which are then deployed via rolling updates or Canary deployments, ensuring zero downtime and allowing for A/B testing of model performance.
    - **Inference Optimization:** Foundry's infrastructure supports optimized inference engines like ONNX Runtime or TensorRT within the Skills model serving containers, maximizing throughput and minimizing latency.
- **5.2. Sandboxed Code Execution & CLI Interaction:**
  - **Foundry's Role:** Provides the isolated and secure environments for Skills to execute arbitrary code (Python, R, shell scripts) or CLI commands.
  - **Microscopic Detail:**
    - **Ephemeral Execution Pods:** When Core instructs a `CodeExecutionSkill` to run a Python script, Foundry's Kubernetes (1.1) spins up a new `Job` or `Pod`.
    - **Docker Containers:** Each execution occurs within a minimal Docker container (e.g., `python:3.10-slim` or a custom image with specific CLI tools).
    - **Linux Namespaces & cgroups:** Docker, running on Foundry's nodes, leverages kernel features like namespaces (PID, Network, Mount, User) for process, network, and filesystem isolation, and cgroups for resource limiting (CPU, memory).
    - **Seccomp-bpf:** Foundry configures a `seccompProfile` in the Kubernetes pod security context, whitelisting only essential system calls that the script requires (e.g., `open`, `read`, `write`, `execve`). Any disallowed syscall (e.g., network access if not explicitly needed, or raw

disk access) will immediately terminate the container. This is a critical layer of defense.

- **Firejail/Bubblewrap (Optional deeper sandboxing):** For extreme isolation requirements, Foundry could configure `Firejail` or `Bubblewrap` to run *inside* the Docker container, adding another layer of user-space sandboxing. This requires careful configuration within the container image.
- **Read-Only Filesystem:** The root filesystem of the execution container is typically mounted as read-only, preventing the script from modifying system binaries or configurations.
- **Temporary Volumes:** Only ephemeral `emptyDir` volumes are provided for temporary file storage, which are destroyed upon pod termination.
- **Network Policy (Egress):** Crucially, the Network Policy (6.2) for these code execution pods typically denies all outbound network connections by default, unless explicitly whitelisted for a specific, verified use case (e.g., downloading a data file from a trusted S3 bucket).
- **Audit Trail:** Every execution event, its duration, resource consumption, and outcome (success/failure) is logged and traced by Foundry's Observability Stack (4.0) and sent to Guardian for analysis.

- **5.3. Data Processing & Analysis Frameworks Integration:**
  - **Foundry's Role:** Provides the scalable compute and data access for Skills to perform large-scale data manipulation and analysis.
  - **Microscopic Detail:**
    - **Distributed Compute:** Foundry can deploy and manage distributed data processing clusters like Apache Spark (via PySpark) or Dask Gateway within Kubernetes. Skills agents requiring distributed computation can submit jobs to these clusters.
    - **Data Access:** Skills agents (whether individual pods or distributed clusters) access data from Foundry's Object Storage (5.1) (e.g., S3-compatible APIs for MinIO or direct cloud APIs). This access is controlled by granular IAM policies (6.6) obtained via Vault (6.1).
    - **Feature Store Integration:** Skills agents can retrieve real-time features from Foundry's Feast Online Store (5.3) for low-latency inference, ensuring consistency with training data.
- **5.4. Generative Task Execution Infrastructure:**
  - **Foundry's Role:** Provides the LLM serving infrastructure for Skills agents that perform generative tasks (e.g., code generation, report writing).
  - **Microscopic Detail:**
    - **Dedicated LLM Services:** Skills might have its own set of specialized LLMs (e.g., CodeLlama, DeepSeek Coder) for code generation. Foundry deploys these using vLLM/Ollama in dedicated Kubernetes deployments with GPU resources, similar to Nexus's LLM but potentially with different model sizes/architectures.

■ **Prompt Engineering Environment:** Foundry provides an environment for Skills developers to test and manage prompt templates for these generative LLMs, ensuring consistent and high-quality output.

---

**6. Foundry ↔ Guardian: The Nexus of Safety, Evaluation & Continuous Learning**

Guardian, Script Kitty's conscience and continuous improvement engine, is inextricably linked with Foundry. Foundry provides the data, compute, and automation necessary for Guardian to fulfill its critical roles.

- **6.1. Policy Enforcement Engine Infrastructure:**
  - **Foundry's Role:** Hosts and provides the underlying infrastructure for Guardian's policy enforcement mechanisms.
  - **Microscopic Detail:**
    - **OPA (Open Policy Agent) Deployment:** Foundry deploys OPA instances as sidecars or dedicated microservices. These OPA instances receive policy rules (Rego code) from a GitOps repository (managed by Foundry's CI/CD).
    - **Policy Decision Points (PDPs):** Critical API endpoints across Script Kitty (e.g., Core's `execute_action` endpoint, Armory's `provision_resource` endpoint) are instrumented to send `ActionRequestSchema` (Protobuf) to Guardian's Policy Enforcement Engine. This engine, running on Foundry, queries its OPA instances (`data.scriptkitty.authz.allow`) for policy decisions.
    - **Real-time Decision:** OPA makes decisions in milliseconds based on the input context (user, action, target resource) and the loaded policies. If a policy is violated, Guardian immediately sends a `DENY` signal back, preventing the action.
- **6.2. Real-time Performance Monitoring & Anomaly Detection Infrastructure:**
  - **Foundry's Role:** Provides the comprehensive observability stack that Guardian consumes for monitoring and anomaly detection.
  - **Microscopic Detail:**
    - **Metrics Consumption:** Guardian directly queries Foundry's Prometheus (4.1) via PromQL to fetch real-time metrics from all Script Kitty components (e.g., `nexus_llm_latency`, `core_planning_failures_total`, `skills_code_execution_errors`).
    - **Log Analysis:** Guardian's anomaly detection models consume enriched logs from Foundry's Loki/Elasticsearch (4.2), looking for patterns indicative of issues (e.g., sudden increase in specific error codes, unusual user behavior patterns).

- - ■ **Trace Analysis:** Guardian can also analyze trace data from Foundry's Jaeger/Zipkin (4.3) to identify performance regressions or unexpected execution paths.
    - ■ **Anomaly Detection Models:** Guardian deploys its own anomaly detection models (e.g., Isolation Forests, LSTM autoencoders) as Skills agents on Foundry's compute (5.1). These models are continuously fed data from Foundry's observability stack.
    - ■ **Alerting Integration:** Guardian can configure its own alerts within Foundry's Prometheus Alertmanager, or consume alerts generated by Foundry and enrich them with further context before escalating.
- **6.3. Automated Evaluation & Training Trigger Infrastructure:**
  - ○ **Foundry's Role:** Provides the MLOps pipelines and compute resources that enable Guardian to trigger and manage model retraining and evaluation.
  - ○ **Microscopic Detail:**
    - ■ **Kubeflow Pipelines/Argo Workflows (ML Orchestration):** Guardian's decision to retrain a model (e.g., based on performance degradation or new data availability) triggers a pre-defined ML pipeline in Foundry's Kubeflow Pipelines or Argo Workflows (5.5).
    - ■ **Pipeline Steps:** A typical pipeline for retraining a Skills agent's model (e.g., a computer vision model):
      - ■ **Data Ingestion:** Reads new training data from Foundry's Object Storage (5.1) (e.g., new images with human labels).
      - ■ **Data Versioning:** Uses DVC (5.2) to `dvc pull` the correct version of the dataset and `dvc commit` any new data, ensuring reproducibility.
      - ■ **Feature Engineering:** Extracts features from the raw data using processing steps that might run on Foundry's Spark/Dask clusters. If a Feature Store (5.3) is used, features are retrieved or computed and pushed to the offline store.
      - ■ **Model Training:** The model training job runs on Foundry's Kubernetes compute, often leveraging GPUs. Parameters are logged to MLflow (5.5).
      - ■ **Hyperparameter Optimization:** Foundry's Ray Tune/Optuna (5.5) can be used to automatically search for optimal hyperparameters during retraining.
      - ■ **Model Evaluation:** The newly trained model is evaluated against a held-out test set (also versioned by DVC) to assess its performance against predefined metrics (accuracy, F1-score, latency). Guardian also monitors for bias using Aequitas/Fairlearn tools on Foundry.
      - ■ **Model Registration:** If the new model performs better and meets safety thresholds, it is registered in a model registry (e.g., MLflow Model Registry).

- ■ **Model Deployment:** The pipeline triggers a new CI/CD build in Foundry (3.0) to package the new model into a Docker image for the relevant Skills inference service, which is then deployed via GitOps.
- ■ **MLflow Integration:** Guardian's training jobs log all parameters, metrics, and model artifacts to Foundry's MLflow (5.5) instance. This allows Guardian (and human operators) to compare different training runs, track experiments, and roll back to previous model versions if needed.
- ● **6.4. Bias Detection & Mitigation Infrastructure:**
  - ○ **Foundry's Role:** Provides the computational environment and access to data for Guardian to run its bias detection tools.
  - ○ **Microscopic Detail:**
    - ■ **Tools Deployment:** Bias detection frameworks like Aequitas, Fairlearn, or IBM's AI Fairness 360 are deployed as specialized containers on Foundry's compute.
    - ■ **Data Access:** These tools access training data and model outputs from Foundry's Data Lake (5.1) and potentially the Feature Store (5.3) to analyze for statistical disparities across different demographic groups.
    - ■ **Reporting:** Guardian generates reports (often visualized in Foundry's Grafana) on detected biases, which can trigger further human review or dedicated retraining efforts aimed at bias mitigation.
- ● **6.5. Human Oversight & Feedback Interface Infrastructure:**
  - ○ **Foundry's Role:** Hosts Guardian's human-in-the-loop interfaces, ensuring secure access and data flow.
  - ○ **Microscopic Detail:**
    - ■ **Web UI Hosting:** Foundry deploys and manages Guardian's web UIs (e.g., built with Streamlit, Plotly Dash, or React/Next.js) within Kubernetes, exposing them via Ingress controllers.
    - ■ **Authentication:** Access to these UIs is secured by Foundry's Keycloak (6.6), enforcing role-based access control (e.g., only "Policy Reviewers" can approve policy overrides).
    - ■ **Data Flow:** The UIs consume data from Foundry's Prometheus (metrics for performance dashboards), Loki/Elasticsearch (logs for detailed debugging), and directly from Guardian's internal APIs (for policy violation reviews, feedback submission).

# Documentation

**Conceptual Code Documentation Focus:**

- **API Schema Definitions:** Detailed Protobuf/JSON schema definitions (`ScriptKittyMessageSchema`, `ScriptKittyOutputSchema`, `ScriptKittyIntentSchema`) for all internal and external message types.
- **Microservice API Contracts:** OpenAPI/Swagger documentation for REST endpoints and `.proto` files for gRPC services exposed by Nexus sub-components.
- **Prompt Engineering Guide:** A living document detailing the system prompts, few-shot examples, and persona instructions for the Nexus LLM. Includes version control for prompts.
- **Configuration Files:** YAML/JSON files for channel adapter credentials, LLM model paths, Redis/PostgreSQL connection strings, rate limits, and security filters.
- **Unit & Integration Test Suites:** Documentation on test coverage for all Nexus functions, including mocked external services.

**Key Technologies, Documentation, and References:**

**1. User Input & Pre-processing Sub-system**

- **1.1. Channel Ingestion Adapters (ChannelAdapter Microservices):**

  - **Description:** Services responsible for translating platform-specific messages into a standardized internal format.

  - **Code Documentation Focus:** Each adapter would have its own repository with:

    - `README.md`: Setup, configuration, and API usage.
    - `schemas/`: Protobuf or Pydantic models for incoming platform-specific payloads and outgoing `ScriptKittyMessageSchema`.
    - `handlers/`: Specific functions for processing webhooks/API events from the platform.
    - `config/`: Environment variables and configuration parameters (e.g., API keys, callback URLs).
    - **Standardized Message Schema (`ScriptKittyMessageSchema`):**
      - **Documentation:** Defined rigorously using Protocol Buffers.
      - **Reference:**
        - [Protocol Buffers Documentation](#)
        - [gRPC Python Tutorial](#)
        - `.proto` file defining `ScriptKittyMessageSchema` fields (`user_id`, `session_id`, `timestamp`, `channel_type`, `raw_text`, `event_type`).

- ○ **Platform-Specific API SDKs:**

  - ■ **Slack API:**
    - ■ **Documentation:** [Slack API Documentation](#)
    - ■ **SDK (Python):** [python-slackclient GitHub](#)
    - ■ **Focus:** `Webhooks`, `Events API`, `OAuth scope configuration`.
  - ■ **Telegram Bot API:**
    - ■ **Documentation:** [Telegram Bot API Documentation](#)
    - ■ **SDK (Python):** [python-telegram-bot GitHub](#)
    - ■ **Focus:** `getUpdates` (long polling) or `setWebhook` (recommended), `sendMessage` method.
  - ■ **FastAPI (for custom Web UI API):**
    - ■ **Documentation:** [FastAPI Documentation](#)
    - ■ **GitHub:** [FastAPI GitHub](#)
    - ■ **Focus:** `Request` and `Response` objects, `Path Operations`, `Dependency Injection`, `Pydantic` for request body validation.
  - ■ **Express.js (for custom Web UI API/Node.js adapters):**
    - ■ **Documentation:** [Express.js Documentation](#)
    - ■ **GitHub:** [Express.js GitHub](#)
    - ■ **Focus:** `middleware`, `routing`, `body-parser` for JSON parsing.
- ● **1.2. Input Normalization & Security Filter:**

  - ○ **Description:** Applies universal pre-processing, data cleaning, and initial security checks.
  - ○ **Code Documentation Focus:**
    - ■ `normalization_rules.py`: Python module detailing `unicode_normalize_text()`, `clean_whitespace()`, `normalize_punctuation()`.
    - ■ `security_filters.py`: Functions for `redact_pii_regex()`, `check_toxicity_heuristic()`, `detect_spam()`.
    - ■ `config/filter_thresholds.yaml`: Configuration for toxicity scores, rate limits.
  - ○ **Open-Source Tools:**
    - ■ **Unidecode (Python):** For ASCII transliterations.
      - ■ **GitHub:** [unidecode GitHub](#)
      - ■ **Focus:** `unidecode.unidecode()` function.
    - ■ **ftfy (fixes text for inconsistent encodings):**
      - ■ **GitHub:** [ftfy GitHub](#)
      - ■ **Focus:** `ftfy.fix_text()` for encoding and formatting issues.

- **Python's built-in string methods:** `str.strip()`, `re.sub()`, `str.lower()`.
  - **Documentation:** [Python String Methods](#)
  - **Documentation:** [Python re module (Regular Expressions)](#)
- **Toxicity/Safety Scoring (Heuristic):**
  - **Example Rule-based:** Custom Python implementation, perhaps using a simple `word_list_matcher`.
  - **Reference:** Concepts from early content moderation heuristics. No specific external tool, but could evolve to use lightweight models like `Hugging Face transformers` if latency allows.

## 2. Dialogue State & Context Management Sub-system

- **2.1. Session Manager:**

  - **Description:** Manages persistent conversational context for each user session.
  - **Code Documentation Focus:**
    - `session_manager_api.proto`: gRPC service definition for `GetSessionContext`, `UpdateSessionContext`, `CreateNewSession`.
    - `models/session_state.py`: Pydantic/Protobuf model for `SessionState` object (including `conversation_history` as list of `ScriptKittyMessageSchema`).
    - `config/database.yaml`: Database connection strings, Redis host/port, TTL settings.
    - `session_manager_db.py`: ORM mappings (SQLAlchemy) for PostgreSQL, Redis client usage.
  - **Open-Source Tools:**
    - **Redis:** For fast, in-memory caching of active sessions.
      - **Documentation:** [Redis Documentation](#)
      - **Client (Python):** [redis-py GitHub](#)
      - **Focus:** `SETEX` (set with expiry), `GET`, `HSET`/`HGET` (hash operations), `JSON.SET`/`JSON.GET` (if RedisJSON module is used for structured data).
    - **PostgreSQL:** For long-term archival and warm storage of older sessions.
      - **Documentation:** [PostgreSQL Documentation](#)
      - **Client (Python ORM):** [SQLAlchemy Documentation](#)
      - **Focus:** Table schema definition (`sessions` table with `session_id` as primary key, `context_data` as JSONB column), indexing for `session_id`, connection pooling.
- **2.2. Contextual Memory & Retrieval:**

- ○ **Description:** Dynamically fetches and updates relevant session data for the LLM.
- ○ **Code Documentation Focus:**
  - ■ `context_builder.py`: Functions for `build_llm_prompt_context()` (e.g., handling sliding window, entity injection).
  - ■ `entity_resolver.py`: Logic for `resolve_entities()`, `update_entity_tracking()`.
- ○ **Open-Source Tools:**
  - ■ **Python's list/dict operations:** For managing conversation history and simple entity tracking within the session object.
  - ■ **Vector Database (Optional/Future):** For advanced long-term memory retrieval.
    - ■ **Chroma:**
      - ■ **Documentation:** [Chroma DB Documentation](Chroma DB Documentation)
      - ■ **GitHub:** [Chroma DB GitHub](Chroma DB GitHub)
      - ■ **Focus:** `Collection API`, `add()`, `query()`, `delete()` for embeddings.
    - ■ **Weaviate:**
      - ■ **Documentation:** [Weaviate Documentation](Weaviate Documentation)
      - ■ **GitHub:** [Weaviate GitHub](Weaviate GitHub)
      - ■ **Focus:** `batching`, `filters`, `nearText` queries.
    - ■ **LlamaIndex / Haystack:** Frameworks for building RAG pipelines.
      - ■ **LlamaIndex:** [LlamaIndex Documentation](LlamaIndex Documentation)
      - ■ **Haystack:** [Haystack Documentation](Haystack Documentation)
      - ■ **Focus:** `Document Loaders`, `Text Splitters`, `Vector Stores`, `Retrievers`.

**3. Core NLU (Natural Language Understanding) & Intent Routing Sub-system**

- ● **3.1. Intent & Entity Extraction (Nexus LLM):**

  - ○ **Description:** Utilizes a fine-tuned LLM for nuanced NLU.
  - ○ **Code Documentation Focus:**
    - ■ `llm_service_api.proto`: gRPC service definition for `InferIntentAndEntities`.
    - ■ `prompts/nexus_nlu_prompt.txt`: The system prompt template with placeholders for session context and user input.
    - ■ `output_schemas/intent_entity_schema.py`: Pydantic/Protobuf model for expected LLM JSON output.
    - ■ `llm_inference_config.yaml`: LLM model ID, temperature, top_p, max_tokens, quantization settings.
  - ○ **Open-Source Tools:**

- **Mistral-7B / Llama 3-8B (or other open-source LLMs):**
  - **Mistral 7B (Mistral AI):** [Mistral AI GitHub](#)
  - **Llama 3 (Meta AI):** [Llama 3 on Hugging Face](#)
  - **Focus:** Model architecture, capabilities, fine-tuning methodologies.
- **vLLM:** For efficient serving of LLMs.
  - **Documentation:** [vLLM Documentation](#)
  - **GitHub:** [vLLM GitHub](#)
  - **Focus:** `Asynchronous API`, `Sampling Parameters`, `Request Batching`, `Serving configuration`.
- **Ollama:** For local LLM serving.
  - **Documentation:** [Ollama Documentation](#)
  - **GitHub:** [Ollama GitHub](#)
  - **Focus:** `REST API`, `Model installation`, `Custom model creation (Modelfile)`.
- **Pydantic:** For validating LLM JSON outputs against a schema.
  - **Documentation:** [Pydantic Documentation](#)
  - **GitHub:** [Pydantic GitHub](#)
  - **Focus:** `BaseModel`, `Field` for schema definition.

- **3.2. Semantic Router & Fallback Mechanism:**

  - **Description:** Intelligently routes queries based on intent confidence.
  - **Code Documentation Focus:**
    - `router_logic.py`: Python module detailing the routing function, confidence thresholds, fallback order.
    - `intent_mapping.yaml`: Mapping of high-confidence intents to internal Nexus handlers or specific Core endpoints.
    - `routing_metrics.py`: Metrics for routing decisions (e.g., `llm_route_count`, `fallback_route_count`, `core_escalation_count`).
  - **Open-Source Tools:**
    - **Hugging Face Transformers (for traditional classifier fallback):**
      - **Documentation:** [Hugging Face Transformers Documentation](#)
      - **GitHub:** [Hugging Face Transformers GitHub](#)
      - **Focus:** `AutoModelForSequenceClassification`, `AutoTokenizer`, `pipeline()` function for text classification.
      - **Pre-trained models:** `bert-base-uncased`, `roberta-base` for fine-tuning on specific intent datasets.
    - **spaCy (for rule-based entity recognition or fallback tokenization):**
      - **Documentation:** [spaCy Documentation](#)
      - **GitHub:** [spaCy GitHub](#)
      - **Focus:** `Matcher`, `EntityRuler`, `Language model loading`.

- **Custom Python/Go code:** The orchestration logic for the hybrid routing.

**4. Response Synthesis & Persona Layer**

- **4.1. Response Generation (Nexus LLM):**

  - **Description:** Generates human-like, persona-consistent responses.
  - **Code Documentation Focus:**
    - `prompts/nexus_response_prompt.txt`: System prompt for generation, incorporating persona guidelines.
    - `response_templates/`: Jinja2 templates for structured responses if needed.
    - `persona_guide.md`: A detailed document outlining Script Kitty's personality traits, communication style, and examples.
  - **Open-Source Tools:**
    - **Mistral-7B / Llama 3-8B:** (Same as for NLU, reusing the serving instance).
    - **Jinja2 (Python):** For flexible templating of responses.
      - **Documentation:** [Jinja2 Documentation](#)
      - **GitHub:** [Jinja2 GitHub](#)
      - **Focus:** `Environment`, `Template.render()`.
- **4.2. Group Chat Facilitation & Multi-Agent Attribution:**

  - **Description:** Creates the "group chat" illusion.
  - **Code Documentation Focus:**
    - `message_formatter.py`: Functions for `add_agent_attribution_prefix()`, `format_progress_update()`.
    - `config/agent_display_names.yaml`: Mapping of internal agent IDs to user-facing display names (e.g., "ScriptKitty.Armory" to "Armory").
  - **Open-Source Tools:** Primarily custom Python/Go logic.

**5. Intermediary & Internal Communication Proxy**

- **5.1. Core Communication Gateway (gRPC Client):**

  - **Description:** Dedicated conduit for all structured communication between Nexus and Script Kitty.Core.
  - **Code Documentation Focus:**
    - `nexus_core_client.py`: Python/Go gRPC client implementation.
    - `proto/`: Directory containing all shared Protobuf definitions (`script_kitty_message.proto`, `task_request.proto`, `task_update.proto`, `task_result.proto`).

- ■ `config/core_grpc_config.yaml`: Core's gRPC server address, port, and security settings (e.g., TLS certificates).
    - ○ **Open-Source Tools:**
        - ■ **gRPC:** High-performance RPC framework.
            - ■ **Documentation:** [gRPC Documentation](#)
            - ■ **Python:** [grpcio GitHub](#)
            - ■ **Focus:** `Channel`, `stub`, `unary_unary`, `stream_stream`, `metadata` for trace context propagation.
        - ■ **Protocol Buffers (Protobuf):** For defining precise schemas.
            - ■ **Documentation:** [Protocol Buffers Documentation](#)
- ● **5.2. User Output Renderer & Channel Dispatcher:**

    - ○ **Description:** Dispatches final responses to the user via the appropriate channel.
    - ○ **Code Documentation Focus:**
        - ■ `output_dispatcher.py`: Function for `dispatch_message_to_channel()`.
        - ■ `schemas/output_schema.py`: Protobuf/Pydantic model for `ScriptKittyOutputSchema`.
        - ■ **Channel Adapter Integrations (Reused):** (See 1.1) Each adapter would have methods like `send_message_to_slack()`, `send_message_to_telegram()`, receiving `ScriptKittyOutputSchema` and translating it.
    - ○ **Open-Source Tools:** Same as for Channel Ingestion Adapters (Slack SDK, Telegram Bot API, FastAPI/Express.js for web sockets).

**Operational Excellence & Scalability (Nexus-Specific):**

- ● **Observability (Integrated with Foundry):**

    - ○ **Prometheus:**
        - ■ **Documentation:** [Prometheus Documentation](#)
        - ■ **Client Libraries (Python):** [prometheus_client GitHub](#)
        - ■ **Focus:** `Counter`, `Gauge`, `Summary`, `Histogram` metrics types, `metrics_middleware` for FastAPI/Express.js.
    - ○ **Grafana:**
        - ■ **Documentation:** [Grafana Documentation](#)
        - ■ **Focus:** Dashboard creation (PromQL queries), Alerting configuration.
    - ○ **Loki:**
        - ■ **Documentation:** [Loki Documentation](#)
        - ■ **GitHub:** [Loki GitHub](#)
        - ■ **Focus:** `LogQL`, labels, `promtail` (log agent).
    - ○ **Fluent Bit / Fluentd:**
        - ■ **Documentation:** [Fluent Bit Documentation](#) / [Fluentd Documentation](#)

- - **GitHub:** [Fluent Bit GitHub](#) / [Fluentd GitHub](#)
    - **Focus:** `Input`, `Filter`, `Output` plugins for log collection.
  - **OpenTelemetry:**
    - **Documentation:** [OpenTelemetry Documentation](#)
    - **SDKs (Python):** [opentelemetry-python GitHub](#)
    - **Focus:** `TracerProvider`, `Span` creation, `Context Propagation`.
  - **Jaeger / Zipkin:**
    - **Jaeger:** [Jaeger Documentation](#)
    - **Zipkin:** [Zipkin Documentation](#)
    - **Focus:** UI for trace visualization, span details.
- **Scalability:**

  - **Kubernetes Horizontal Pod Autoscaler (HPA):**
    - **Documentation:** [Kubernetes HPA Documentation](#)
    - **Focus:** `apiVersion: autoscaling/v2`, `minReplicas`, `maxReplicas`, `metrics` (CPU, memory, custom metrics via `metrics.k8s.io` API).
  - **Kubernetes Deployments:**
    - **Documentation:** [Kubernetes Deployments Documentation](#)
    - **Focus:** `replicas`, `selector`, `template`, `strategy` (RollingUpdate).
  - **Service Mesh (Istio/Linkerd):** For advanced load balancing.
    - **Istio:** [Istio Documentation](#)
    - **Linkerd:** [Linkerd Documentation](#)
    - **Focus:** `VirtualService`, `DestinationRule` for traffic management, `ServiceEntry` for external services.
- **Continuous Evolution (Via Script Kitty.Foundry):**

  - **Docker:**
    - **Documentation:** [Docker Documentation](#)
    - **Focus:** `Dockerfile` best practices, `docker build`, `docker push`.
  - **CI/CD Pipelines (e.g., GitLab CI/CD, GitHub Actions):**
    - **GitLab CI/CD:** [GitLab CI/CD Documentation](#)
    - **GitHub Actions:** [GitHub Actions Documentation](#)
    - **Focus:** `workflows` (GitHub Actions), `jobs`, `stages`, `steps`, `triggers`, `artifacts`.
  - **GitOps (Argo CD / Flux CD):**
    - **Argo CD:** [Argo CD Documentation](#)
    - **Flux CD:** [Flux CD Documentation](#)
    - **Focus:** `Application` resource (Argo CD), `HelmRelease` (Flux CD), `reconciliation loops`, `sync policies`, `rollbacks`.

- **Model Management (part of Foundry's MLOps):** Nexus LLM model updates are treated as new container images deployed via the Foundry's CI/CD and GitOps.

# File tree

**Overall Architectural Philosophy: The Monorepo with Strategic Boundaries**

For Script Kitty, a **monorepo** strategy is meticulously chosen as the primary repository layout. This decision is not arbitrary; it's a deliberate choice to foster:

1. **Atomic Changes:** Facilitates large-scale refactorings or feature implementations that span multiple services, ensuring consistency across Nexus, Core, Armory, Skills, Guardian, and Foundry. A single commit can update shared Protobuf definitions and the services that consume them.
2. **Simplified Dependency Management:** Common libraries, Protocol Buffer definitions, and shared utility functions can be centrally managed and easily consumed by all components.
3. **Unified CI/CD & MLOps:** Enables a holistic view and orchestration of build, test, and deployment pipelines across the entire AGI, managed by Foundry. New models or policy updates can trigger integrated validation across all affected components.
4. **Enhanced Discoverability & Collaboration:** Developers can easily navigate the entire codebase, understand inter-service communication, and contribute across different domains.
5. **Centralized Security Auditing:** A single repository simplifies security scans and policy enforcement across the entire codebase.

However, recognizing the sheer scale and distinct deployment lifecycles of individual microservices, the monorepo will be structured with **clear, isolated build contexts** for each major component. This allows for independent Docker image builds and targeted deployments, preserving the benefits of microservices while leveraging the monorepo's advantages.

---

## Top-Level Repository Layout: The Script Kitty Root

The root of the `script-kitty-agi/` repository is the nexus of its creation, providing immediate insight into the modular architecture.

```
/script-kitty-agi
├── .git/                    # Git version control metadata
├── .github/                   # GitHub-specific configurations (e.g., issue templates, pull request
templates)
│   ├── workflows/            # GitHub Actions CI/CD workflows (if GitHub is the chosen Git
platform)
│   │   ├── ci_build.yaml
│   │   ├── cd_deploy_prod.yaml
│   │   └── security_scan.yaml
├── .gitlab/                 # GitLab-specific configurations (if GitLab is the chosen Git platform)
│   ├── ci/                  # GitLab CI/CD includes
```

```
│   │   ├── build.gitlab-ci.yml
│   │   └── deploy.gitlab-ci.yml
│   ├── issue_templates/
│   └── merge_request_templates/
├── .devcontainer/          # DevContainer configuration for standardized development
environments
│   ├── devcontainer.json
│   └── Dockerfile
├── .husky/                 # Git hooks (e.g., pre-commit linting, pre-push tests)
├── .vscode/                # Visual Studio Code workspace settings and extensions
recommendations
├── docs/                   # Comprehensive system documentation (user guides, architecture,
API specs, MLOps playbooks)
│   ├── architecture/       # High-level architecture diagrams, component descriptions
│   │   ├── overview.md
│   │   └── component_flow.excalidraw
│   ├── api/                # API documentation (OpenAPI/Swagger specs for external APIs)
│   │   ├── nexus_api.yaml
│   │   └── guardian_policy_api.yaml
│   ├── deployment/         # Deployment guides, Helm chart usage, Kubernetes cluster
setup
│   │   └── production_deployment.md
│   ├── development/        # Contribution guidelines, local setup, testing methodologies
│   │   └── contributing.md
│   ├── mloppl_ops/         # MLOps best practices, model lifecycle management
│   │   └── model_lifecycle.md
│   └── README.md           # Root README, project overview, quick start
├── third_party/            # Vendored third-party libraries or external dependencies not
managed by package managers
├── tools/                  # Project-specific helper scripts (e.g., proto generation, data
migration tools, local dev setups)
│   ├── proto_gen.sh
│   ├── apply_k8s_manifests.sh
│   └── local_dev_setup.sh
├── .editorconfig           # Consistent coding styles across different IDEs/editors
├── .gitignore              # Files/directories to be ignored by Git
├── Makefile                # Top-level Makefile for common commands (build all, test all,
clean)
├── README.md               # Primary project README, overall vision, high-level
architecture
├── LICENSE                 # Project license (e.g., Apache 2.0, MIT)
├── CODE_OF_CONDUCT.md      # Guidelines for community behavior
├── CONTRIBUTING.md         # Detailed contribution guidelines
└── components/             # The heart of Script Kitty: individual microservice components
```

```
├── nexus/               # Script Kitty.Nexus component
├── core/                # Script Kitty.Core component
├── armory/               # Script Kitty.Armory component
├── skills/            # Script Kitty.Skills component
├── guardian/             # Script Kitty.Guardian component
└── foundry/               # Script Kitty.Foundry component (core MLOps services)
```

**Microscopic Dissection of Top-Level Directories:**

- **`.git/`**: The `.git` directory holds all the information about the Git repository, including objects (actual data), refs (pointers to commits), and configuration. It's the immutable record of Script Kitty's evolution, allowing for precise version control, history tracking, and collaborative development.
- **`.github/` / `.gitlab/`**: These directories contain platform-specific CI/CD workflow definitions.
    - `workflows/`: Each `.yaml` file defines a GitHub Actions workflow.
        - `ci_build.yaml`: Triggers on push/PR, runs linting, unit/integration tests, and builds Docker images for all changed components, pushing them to Foundry's secure container registry.
        - `cd_deploy_prod.yaml`: Triggers on merged main branch, or manual approval, initiating production deployment via GitOps agents (Argo CD/Flux CD) managed by Foundry. This pipeline will interact with the `kubernetes/` manifests.
        - `security_scan.yaml`: Scheduled nightly or on push, runs static application security testing (SAST), dependency vulnerability scans (via Trivy/Clair), and potentially code secret scanning across the entire monorepo.
    - These CI/CD configurations are the very instructions for Foundry's automated pipelines, ensuring every code change is rigorously validated and deployed with precision.
- **`.devcontainer/`**: Essential for providing a standardized development environment.
    - `devcontainer.json`: Defines the Docker image, VS Code extensions, and port forwards for a consistent development experience across all contributors, ensuring everyone works with the same tool versions and dependencies.
    - `Dockerfile`: Builds the custom development container image, pre-installing all necessary SDKs, compilers, and tools (e.g., Python, Go, Node.js, Protobuf compiler, Kubernetes CLI tools, DVC, Helm, Terragrunt, specific LLM quantization tools).
- **`.husky/`**: Integrates Git hooks (e.g., `pre-commit`, `pre-push`) to enforce quality gates *before* code even reaches the CI/CD pipeline.

- ○ `pre-commit`: Runs linters (`flake8`, `mypy`), formatters (`black`), and basic unit tests on staged changes, ensuring code quality standards are met at the developer's workstation.
- ○ `pre-push`: Runs more extensive tests or pushes validated code to a temporary branch, preventing broken code from entering the main branch.
- **`.vscode/`**: Contains workspace-specific settings for Visual Studio Code, ensuring consistency in formatting, linting, and debugger configurations across the development team.
- **`docs/`**: The single source of truth for all documentation.
  - ○ `architecture/`: Markdown and diagram files detailing Script Kitty's layered architecture, communication flows (e.g., sequence diagrams for `Nexus -> Core -> Skills` interactions), and design patterns. Crucial for onboarding new developers and maintaining long-term understanding.
  - ○ `api/`: OpenAPI/Swagger specifications for any public or internal APIs exposed by Script Kitty. Nexus's user-facing API, Guardian's policy management API, or Armory's external tool API definitions would reside here, ensuring consistent API contracts.
  - ○ `deployment/`: Detailed guides for deploying Script Kitty to various environments (local, staging, production) using Foundry's Kubernetes manifests and Helm charts. Includes prerequisites, step-by-step instructions, and troubleshooting.
  - ○ `development/`: Guidelines for contributing code, setting up local development environments, running tests, and adhering to coding standards.
  - ○ `mlops_ops/`: Best practices and runbooks for MLOps operations: model retraining, A/B testing, feature engineering workflows, monitoring model drift, and handling data biases.
  - ○ `README.md`: Provides a high-level overview of the entire Script Kitty project, its mission, core components, and quick links to more detailed documentation.
- **`third_party/`**: For vendored dependencies not easily managed by package managers, or for specialized binaries. This could include pre-compiled binaries for specific hardware (e.g., custom GPU drivers, quantum computing SDK components) or patched libraries.
- **`tools/`**: A collection of cross-component helper scripts.
  - ○ `proto_gen.sh`: A script that orchestrates the generation of client/server code from `.proto` files (located in `shared/protobuf/`) for all relevant languages (Python, Go). This ensures communication schemas are always in sync.
  - ○ `apply_k8s_manifests.sh`: A helper for applying Kubernetes YAMLs in a specific order during local development or for quick testing.
  - ○ `local_dev_setup.sh`: Script to automate the setup of local dependencies (e.g., Minikube, kind, Redis, PostgreSQL) for quick local development of Script Kitty.

- **`.editorconfig`**: Ensures consistent coding styles (indentation, line endings, character sets) across different editors and IDEs used by contributors.
- **`.gitignore`**: Specifies intentionally untracked files that Git should ignore (e.g., build artifacts, temporary files, environment-specific configurations, Python `__pycache__`).
- **`Makefile`**: A top-level orchestrator for common tasks across the entire monorepo (e.g., `make build-all`, `make test-all`, `make clean`, `make deploy-local`).
- **`README.md` (root)**: The primary entry point for anyone discovering the Script Kitty project, providing an overview of its vision, goals, and high-level architecture.
- **`LICENSE` / `CODE_OF_CONDUCT.md` / `CONTRIBUTING.md`**: Standard open-source project files, defining legal terms, community behavior, and contribution guidelines.
- **`components/`**: The core directory housing all individual microservices of Script Kitty. Each sub-directory here represents a distinct, deployable component.

---

## Detailed Component-Specific Structures: Inside `components/`

Each sub-directory within `components/` follows a largely consistent, yet specialized, microservice pattern. This consistency simplifies development, deployment, and operational management.

```
/script-kitty-agi
└── components/
    ├── nexus/
    │   ├── src/
    │   │   ├── api/               # REST API endpoints (e.g., for web chat, external integrations)
    │   │   │   └── v1/
    │   │   │       └── router.py
    │   │   ├── channel_adapters/     # Adapters for various messaging platforms
    │   │   │   ├── slack_adapter.py
    │   │   │   ├── web_socket_adapter.py
    │   │   │   └── telegram_adapter.py
    │   │   ├── nlu/               # Natural Language Understanding logic
    │   │   │   ├── intent_router.py   # LLM-based and traditional NLU routing
    │   │   │   └── entity_extractor.py
    │   │   ├── conversation_manager/  # Dialogue state and context management
    │   │   │   ├── session_manager.py
    │   │   │   └── context_store.py
    │   │   ├── response_generator/    # Persona and response synthesis
    │   │   │   ├── persona_handler.py
    │   │   │   └── response_synthesizer.py
    │   │   ├── llm_client/         # Client for Nexus LLM inference service
    │   │   │   └── vllm_client.py
```

```
│   │   ├── main.py              # Entry point for the Nexus service
│   │   └── config.py            # Configuration parsing and loading
│   ├── tests/
│   │   ├── unit/
│   │   ├── integration/
│   │   └── e2e/                 # End-to-end tests for Nexus
│   ├── deploy/                  # Kubernetes manifests and Helm charts for Nexus deployment
│   │   ├── k8s/
│   │   │   ├── deployment.yaml
│   │   │   ├── service.yaml
│   │   │   └── hpa.yaml
│   │   └── helm/
│   │       ├── Chart.yaml
│   │       ├── values.yaml
│   │       └── templates/
│   ├── Dockerfile               # Docker build context for Nexus
│   ├── requirements.txt         # Python dependencies for Nexus
│   ├── pyproject.toml           # Poetry/PDM configuration
│   └── README.md                # Nexus-specific README
│
├── core/
│   ├── src/
│   │   ├── planning_engine/     # HTN planning and goal decomposition
│   │   │   ├── htn_planner.py
│   │   │   └── plan_templates.py
│   │   ├── orchestration_manager/ # Agent lifecycle and task allocation
│   │   │   └── agent_orchestrator.py
│   │   ├── knowledge_graph_client/ # GKGS interaction
│   │   │   └── gkgs_client.py
│   │   ├── internal_api/        # gRPC server for inter-service communication
│   │   │   └── core_service.proto # Protobuf definition for Core's internal API
│   │   ├── task_executor/       # Dispatches tasks to Armory/Skills/Guardian
│   │   │   └── task_dispatcher.py
│   │   ├── feedback_integrator/   # Processes feedback from Guardian
│   │   │   └── feedback_handler.py
│   │   ├── main.py
│   │   └── config.py
│   ├── models/                  # Core's internal LLM/planning models (e.g., Llama 3-70B
weights if self-hosted)
│   ├── tests/
│   ├── deploy/
│   ├── Dockerfile
│   ├── requirements.txt
│   └── README.md
```

```
├── armory/
│   ├── src/
│   │   ├── web_scraper/          # Intelligent web scraping and parsing
│   │   │   ├── browser_automator.py
│   │   │   └── content_parser.py
│   │   ├── data_digestion/       # NLP for data extraction and summarization
│   │   │   ├── entity_extractor.py
│   │   │   └── summarizer.py
│   │   ├── tool_registry/        # Dynamic tool catalog management
│   │   │   ├── tool_db.py
│   │   │   └── manifest_schema.py # Schema for ToolManifest
│   │   ├── wrapper_generator/    # LLM-based code generation for wrappers
│   │   │   ├── code_templates/   # Secure wrapper templates
│   │   │   │   └── python_api_wrapper.j2
│   │   │   └── llm_codegen.py
│   │   ├── resource_provisioner/ # IaC orchestration
│   │   │   ├── terraform_client.py
│   │   │   └── ansible_client.py
│   │   ├── internal_api/         # gRPC service for Core to interact with Armory
│   │   │   └── armory_service.proto
│   │   ├── main.py
│   │   └── config.py
│   ├── data/                     # Sample data for testing scraping/digestion
│   ├── tools_manifests/          # Declarative definitions of known tools (seed data)
│   │   └── cli_tool_example.yaml
│   ├── tests/
│   ├── deploy/
│   ├── Dockerfile
│   ├── requirements.txt
│   └── README.md
│
├── skills/
│   ├── src/
│   │   ├── model_inference/      # General model serving framework
│   │   │   └── infer_engine.py
│   │   ├── sandboxed_execution/  # Code execution environment
│   │   │   ├── docker_executor.py
│   │   │   └── safe_evaluator.py
│   │   ├── data_analysis/        # Libraries for data processing (Pandas, NumPy, etc.)
│   │   │   └── data_processor.py
│   │   ├── generative_tasks/     # LLM for code/report generation
│   │   │   └── code_generator.py
│   │   ├── internal_api/         # gRPC service for Core to interact with Skills
```

```
│   │   │   └── skills_service.proto
│   │   ├── main.py
│   │   └── config.py
│   ├── models/                # Specific ML models (e.g., CV models, NLP embeddings,
fine-tuned code LLMs)
│   │   ├── computer_vision_model_v1.pkl
│   │   └── text_embedding_model_v2.pt
│   ├── tests/
│   ├── deploy/
│   ├── Dockerfile             # Includes dependencies for sandboxing (e.g., Firejail)
│   ├── requirements.txt
│   └── README.md
│
├── guardian/
│   ├── src/
│   │   ├── policy_engine/      # OPA integration and custom rule engine
│   │   │   ├── opa_client.py
│   │   │   └── policy_rules.rego  # OPA policy definitions
│   │   ├── monitoring_processor/  # Real-time metric/log anomaly detection
│   │   │   ├── anomaly_detector.py
│   │   │   └── metrics_aggregator.py
│   │   ├── evaluation_manager/   # Automated model evaluation and trigger logic
│   │   │   ├── model_evaluator.py
│   │   │   └── retraining_trigger.py
│   │   ├── feedback_interface/   # Human oversight and feedback processing
│   │   │   └── feedback_handler.py
│   │   ├── bias_detection/      # Tools for bias analysis in data/models
│   │   │   └── bias_analyzer.py
│   │   ├── internal_api/        # gRPC for Core/other agents to request policy checks
│   │   │   └── guardian_service.proto
│   │   ├── main.py
│   │   └── config.py
│   ├── policies/              # Externalized policy definitions (Rego files, YAML rules)
│   │   ├── ethical_guidelines.rego
│   │   └── data_privacy_rules.yaml
│   ├── dashboards/            # Grafana dashboard definitions (JSON) for Guardian's
metrics
│   │   └── guardian_overview.json
│   ├── tests/
│   ├── deploy/
│   ├── Dockerfile
│   ├── requirements.txt
│   └── README.md
│
│
```

```
└── foundry/
    ├── src/
    │   ├── kubernetes_manager/    # API clients for Kubernetes orchestration
    │   │   └── k8s_client.py
    │   ├── ci_cd_orchestrator/    # Triggers and monitors CI/CD pipelines
    │   │   └── pipeline_trigger.py
    │   ├── observability_agent/   # Prometheus/Loki/OpenTelemetry agents
    │   │   ├── metrics_exporter.py
    │   │   └── log_collector.py
    │   ├── data_manager/          # Interfaces with data lake, DVC, Feature Store
    │   │   ├── object_storage_client.py
    │   │   └── dvc_integration.py
    │   ├── security_manager/      # Integrates with Vault, Keycloak, network policies
    │   │   ├── vault_client.py
    │   │   └── keycloak_client.py
    │   ├── ml_pipeline_orchestrator/ # Kubeflow/Argo Workflows client
    │   │   └── argo_workflows_client.py
    │   ├── main.py
    │   └── config.py
    ├── kubernetes/                # Core Kubernetes manifests for the entire Script Kitty system
    │   ├── base/                  # Baseline cluster setup (e.g., cert-manager, nginx-ingress)
    │   │   ├── namespace.yaml
    │   │   └── network_policies/  # Global network policies for all services
    │   │       └── default_deny.yaml
    │   ├── argo_cd/               # Argo CD deployment manifests
    │   ├── prometheus_grafana/    # Monitoring stack deployment
    │   ├── loki_fluentbit/        # Logging stack deployment
    │   ├── jaeger/                # Tracing backend deployment
    │   ├── vault/                 # HashiCorp Vault deployment
    │   ├── keycloak/              # Keycloak identity management deployment
    │   ├── object_storage/        # MinIO/Ceph deployment (if self-hosted)
    │   ├── feature_store/         # Feast deployment
    │   ├── ml_pipelines/          # Kubeflow/Argo Workflows deployment
    │   ├── common/                # Shared resources like service accounts, cluster roles
    │   │   └── service_accounts.yaml
    │   └── crds/                  # Custom Resource Definitions for Script Kitty (e.g.,
ScriptKittyAgent)
    │       └── scriptkitty_agent_crd.yaml
    ├── helm_charts/               # Helm charts for deploying the entire Script Kitty system
(meta-chart)
    │   ├── script-kitty-agi/
    │   │   ├── Chart.yaml
    │   │   ├── values.yaml
    │   │   └── templates/
```

```
│  │      ├── _helpers.tpl
│  │      ├── nexus_subchart.yaml
│  │      └── core_subchart.yaml
│  └── subcharts/          # Individual component helm charts (optional, could be in
component/deploy)
    ├── data/              # Example data for Foundry (e.g., for testing data pipelines)
    ├── tests/
    ├── Dockerfile
    ├── requirements.txt
    └── README.md
```

**Microscopic Dissection of Component Directories (`components/<component_name>/`):**

Each component directory (e.g., `nexus/`, `core/`) is a self-contained unit designed for independent development, testing, and deployment.

- **`src/`**: The core source code for the component.
    - **Sub-directories within `src/`**: These align with the functional sub-systems detailed in previous analyses. For example, `nexus/src/nlu/` would contain the LLM inference client (`llm_client.py`) and the routing logic (`intent_router.py`). The granular separation reflects the single responsibility principle for microservices.
    - **`main.py`**: The primary entry point for the service, typically initializing configurations, setting up logging, registering API routes/gRPC services, and starting the main event loop.
    - **`config.py`**: Handles loading configuration parameters from environment variables, Kubernetes ConfigMaps, or secrets (accessed securely via Foundry's Vault client). This ensures environment-specific settings are injected at runtime.
    - **`internal_api/`**: Contains the `.proto` files defining the gRPC service contracts for inter-component communication. These are the explicit, versioned APIs that Core uses to interact with Armory, Skills, and Guardian, and that Nexus uses to talk to Core. The `proto_gen.sh` script (from `tools/`) compiles these into language-specific client/server stubs.
- **`models/`**: Specific to components that manage or host ML models.
    - **`core/models/`**: Could contain the weights for Core's planning LLM (if self-hosted) or references to where they are pulled from Foundry's object storage.
    - **`skills/models/`**: Stores specific ML models (e.g., computer vision models, NLP models, fine-tuned code generation models) that a particular Skills agent serves. These would typically be versioned and loaded via Foundry's data management features.
- **`data/`**: Small, component-specific datasets or sample files primarily for local testing.

- ○ `armory/data/`: Sample web pages for testing scraping, or small datasets for digestion tests.
- ● **`tests/`**: A critical directory ensuring code quality and correctness.
  - ○ **`unit/`**: Isolated tests for individual functions and classes, ensuring their logic is sound.
  - ○ **`integration/`**: Tests interactions between components within the same service (e.g., Nexus's NLU communicating with its Session Manager).
  - ○ **`e2e/` (End-to-End)**: High-level tests that simulate real user flows across multiple components (e.g., a user query to Nexus, processed by Core, executed by Skills, and the response returned), ensuring the entire chain works as expected. These tests are vital for AGI-level validation.
- ● **`deploy/`**: Contains the Kubernetes manifests and Helm charts specific to deploying *this individual component*.
  - ○ **`k8s/`**: Raw Kubernetes YAML definitions (`deployment.yaml`, `service.yaml`, `hpa.yaml`, `ingress.yaml`). These are the declarative blueprints for how Kubernetes should run the service.
  - ○ **`helm/`**: A self-contained Helm chart for the component.
    - ■ `Chart.yaml`: Defines the chart metadata.
    - ■ `values.yaml`: Default configurable parameters for the component's deployment (e.g., replica count, resource limits, image tag).
    - ■ `templates/`: Contains Jinja2-templated Kubernetes YAMLs that consume values from `values.yaml` and `.helmignore`.
- ● **`Dockerfile`**: Defines the Docker image for the component. This is the precise instruction set for containerizing the service, specifying base image, dependencies, environment variables, and the entry point. Each `Dockerfile` is optimized for size and security (e.g., multi-stage builds).
- ● **`requirements.txt` / `pyproject.toml`**: Specifies the exact Python dependencies for the component. This ensures reproducible builds and environments, crucial for preventing dependency conflicts in a monorepo.
- ● **`README.md`**: Provides a component-specific README, detailing its purpose, APIs, configuration, and how to run it locally.

---

## Shared Resources & Infrastructure: The Glue of Script Kitty

Beyond individual components, Script Kitty relies on a set of shared resources and infrastructure configurations that bind the system together.

```
/script-kitty-agi
├── shared/
│   ├── protobuf/          # All Protobuf definitions for inter-service communication
```

```
│   │   ├── script_kitty_messages.proto  # Core message schemas
(ScriptKittyMessageSchema, TaskRequestSchema, etc.)
│   │   ├── nexus_service.proto
│   │   ├── core_service.proto
│   │   ├── armory_service.proto
│   │   ├── skills_service.proto
│   │   ├── guardian_service.proto
│   │   └── security_common.proto  # Common security-related enums/messages
│   ├── python_libs/          # Shared Python libraries/utilities (e.g., logging setup, common
validators)
│   │   ├── sk_logging.py
│   │   └── sk_exceptions.py
│   ├── go_libs/              # Shared Go libraries/utilities (if Go is used for some services)
│   ├── common_types/          # Language-agnostic type definitions (e.g., shared
enums/constants)
│   │   └── task_status_enum.yaml
│   └── templates/            # Shared templating resources (e.g., for LLM prompting)
│       └── base_persona_prompt.j2
├── kubernetes/              # Global Kubernetes infrastructure manifests (Managed by
Foundry)
│   ├── base/                # Core cluster setup
│   ├── argo_cd/             # GitOps agent deployment
│   ├── prometheus_grafana/      # Observability stack
│   ├── loki_fluentbit/
│   ├── jaeger/
│   ├── vault/               # Secrets management
│   ├── keycloak/            # Identity management
│   ├── object_storage/        # Self-hosted object storage (MinIO)
│   ├── feature_store/         # Feast deployment
│   ├── ml_pipelines/          # Kubeflow/Argo Workflows deployment
│   ├── common/              # Shared Kubernetes resources
│   └── crds/              # Script Kitty's custom CRDs
└── data_store/             # High-level directory for data artifacts (Data Lake, DVC cache,
Feature Store data)
    ├── raw_data/            # Raw ingested data (e.g., web scrapes, user inputs)
    ├── processed_data/        # Cleaned, transformed datasets
    ├── features/            # Feature store data (offline)
    ├── model_artifacts/        # Trained model checkpoints, serialized models
    ├── dvc_cache/           # DVC's internal cache for data files
    └── experiment_logs/        # MLflow/DVC experiment tracking logs
```

**Microscopic Dissection of Shared Resources:**

- **shared/**: The nexus for inter-component collaboration and consistency.
  - **protobuf/**: This is the absolute cornerstone of inter-service communication within Script Kitty. Every `.proto` file here defines a gRPC service and its message types (`ScriptKittyMessageSchema`, `TaskRequestSchema`, `ToolManifestSchema`). Foundry's CI/CD pipeline runs `proto_gen.sh` automatically on changes, ensuring all language-specific client/server stubs are re-generated and components rebuild if their interfaces change. This guarantees strict type-checking and robust, efficient communication.
  - **python_libs/** / **go_libs/**: Centralized location for common utility functions, logging configurations, custom exceptions, and base classes that are shared across multiple services to prevent code duplication and enforce consistency.
  - **common_types/**: Language-agnostic definitions (e.g., YAML files for enums, JSON schemas for configuration blocks) that can be parsed by different services.
  - **templates/**: Shared templates for LLM prompting (e.g., base persona, common instruction formats) or code generation, ensuring consistent output across components.
- **kubernetes/**: The declarative definition of Script Kitty's entire Kubernetes infrastructure. This directory is primarily managed by Foundry's GitOps agent (Argo CD/Flux CD).
  - **base/**: Contains manifests for foundational services (e.g., `namespace.yaml` for `script-kitty-prod`, `network_policies/default_deny.yaml` to establish a zero-trust model by default).
  - **argo_cd/**, **prometheus_grafana/**, **loki_fluentbit/**, **jaeger/**, **vault/**, **keycloak/**, **object_storage/**, **feature_store/**, **ml_pipelines/**: These sub-directories contain the detailed Kubernetes manifests (Deployments, StatefulSets, Services, Ingresses, ConfigMaps, Secrets, PVCs) for deploying each foundational MLOps tool that Foundry leverages. Each manifest specifies resource requests/limits, health checks, readiness probes, and security contexts to ensure robust operation.
  - **common/**: Kubernetes resources shared across multiple Script Kitty services, such as `ServiceAccount` definitions, `ClusterRole` and `ClusterRoleBinding` for RBAC, ensuring that components have exactly the permissions they need and no more.
  - **crds/**: The Custom Resource Definitions for Script Kitty's high-level abstractions, such as `scriptkitty_agent_crd.yaml`, which allows Core to declare desired agent types and have Kubernetes provision them via custom controllers.
- **data_store/**: A logical grouping for the data lake components and artifacts. While the actual data resides in object storage (MinIO, S3, etc.), this directory represents how the data is *structured and versioned* within the project's operational paradigm.

- ○ **`raw_data/`**: Ingested, untouched data. This might contain data from Armory's web scraping or initial user inputs, preserved in its original form.
- ○ **`processed_data/`**: Cleaned, transformed, and aggregated datasets ready for feature engineering or direct model training.
- ○ **`features/`**: Contains the definitions and (offline) materialized views of features managed by Feast.
- ○ **`model_artifacts/`**: Stores versioned model checkpoints and serialized models from training runs orchestrated by Guardian. This is where models are registered for deployment by Foundry.
- ○ **`dvc_cache/`**: The internal cache directory for DVC, which stores the actual data files linked by `.dvc` files. This directory is typically large and often placed on a separate persistent volume.
- ○ **`experiment_logs/`**: Outputs from MLflow or DVC experiment tracking, including run metadata, metrics, and parameters, crucial for Guardian to analyze model performance and decide on retraining.

---

## Flawless Execution through Meticulous Structure

This microscopic dissection of Script Kitty's file tree and project structure reveals a deliberate design choice that underpins its ambitious capabilities:

- **Modularity & Decoupling:** Each component is a distinct, deployable unit, promoting independent development and scaling while adhering to clear API contracts defined in `shared/protobuf/`.
- **Reproducibility:** From standardized development containers (`.devcontainer/`) to versioned data (`data_store/`, DVC), and immutable Docker images, every step of Script Kitty's lifecycle is designed to be fully reproducible. This is critical for debugging complex AGI behaviors and for auditing.
- **Automation-First:** The extensive use of CI/CD (`.github`/`.gitlab/`), GitOps (`kubernetes/` manifests managed by Argo CD), and internal tooling (`tools/`) ensures that human error is minimized, and all changes are validated and deployed automatically by Foundry.
- **Observability & Debuggability:** The structured logging, metrics, and tracing (`observability_stack/`) are woven into every component, providing unparalleled visibility into Script Kitty's internal workings, allowing rapid identification and resolution of any issues that may arise in its intelligent processes.
- **Scalability:** The microservice architecture, combined with Kubernetes HPA/VPA and Foundry's robust infrastructure orchestration, allows each component to scale independently based on demand, optimizing resource utilization.
- **Security by Design:** Layered security controls from the file system (`.git/ignore`), through code analysis (CI/CD scans), runtime protection (Falco), network policies

(Calico/Cilium), and centralized secrets/identity management (Vault/Keycloak), are deeply embedded in the structure.

- **Continuous Learning & Evolution:** The explicit structure for `data_store/`, `models/`, and the integration with MLOps tools (Feast, DVC, Kubeflow/Argo Workflows) directly supports Guardian's ability to trigger and manage continuous retraining and improvement of Script Kitty's intelligent capabilities.

This granular blueprint is not just a theoretical exercise; it is the practical manifestation of how Script Kitty will be built and operated, ensuring its journey towards an open-source AGI is executed with surgical precision and unwavering reliability. This structure is the foundational genius, enabling the seamless, evolving, and supremely capable Script Kitty we envision.

# Dependencies

## Conceptual Model:

- **Requirements (Rx):** The explicit and implicit needs of a component or function to perform its designated role. These can be computational (CPU, GPU, RAM), data-related (specific input schemas, access to datasets), operational (latency, throughput), or environmental (network access, security context).
- **Dependencies (Dx):** The specific components, services, or data sources that a given component or function *relies upon* to fulfill its requirements. These can be:
  - **Internal Dependencies (DIx):** Other Script Kitty components (Nexus, Core, Armory, Skills, Guardian, Foundry services).
  - **External Dependencies (DEx):** Third-party APIs, external databases, cloud services, open-source libraries, or proprietary tools.
- **Correlations (Cx):** The dynamic relationships, feedback loops, and cascading effects between components. How the output, state, or performance of one component directly influences the requirements, behavior, or even the evolution of another. This is where the true "intelligence" of the system's integration manifests.

---

## Overall System-Level Core Requirements & Dependencies:

Before diving into individual components, it's crucial to acknowledge the overarching requirements and the fundamental dependencies that permeate the entire Script Kitty ecosystem:

- **Universal Requirement:** Highly Available, Scalable, Secure, and Performant Compute Infrastructure.
  - **Correlation:** Directly provided by **Script Kitty.Foundry**'s Kubernetes-Native Compute & Orchestration, Container Registry, CI/CD, and Security services. Without Foundry, no other component could even exist or operate reliably.
- **Universal Requirement:** Robust Inter-Service Communication.
  - **Correlation:** Primarily provided by **Script Kitty.Foundry**'s Kubernetes networking, Service Mesh (Istio/Linkerd), and universally enforced **gRPC/Protobuf** standards. Every component relies on this for reliable communication with every other component.
- **Universal Requirement:** Comprehensive Observability (Metrics, Logging, Tracing).
  - **Correlation:** Directly provided by **Script Kitty.Foundry**'s Observability Stack. All components *emit* data according to Foundry's standards, and Foundry *consumes, aggregates, and visualizes* it. This is a critical feedback loop for **Script Kitty.Guardian**'s monitoring and evaluation.
- **Universal Requirement:** Secure Credential Management.
  - **Correlation:** Met by **Script Kitty.Foundry**'s Security & Identity Management (HashiCorp Vault, Keycloak). Every component that needs to access external

APIs or sensitive internal resources depends on Vault for secrets and Keycloak for identity.

- **Universal Requirement:** Continuous Improvement & Learning Mechanism.
  - **Correlation:** Orchestrated by **Script Kitty.Core**'s feedback integration, driven by **Script Kitty.Guardian**'s evaluation, and underpinned by **Script Kitty.Foundry**'s MLOps pipelines (Data Management, CI/CD, ML Pipelines). This is the meta-requirement for Script Kitty to be an "AGI."

---

**Detailed Dissection by Component:**

# I. Script Kitty.Nexus - The User's Portal & Group Chat Facilitator

**Core Purpose (Recap):** The seamless, unified conversational interface, abstracting multi-agent complexity.

**1. User Input & Pre-processing Sub-system:**

- **1.1. Channel Ingestion Adapters (ChannelAdapter Microservices):**

  - **Requirements (Rx):**
    - **Functional:** Real-time reception of messages from diverse platforms (Slack, Web, Telegram, REST API). Conversion to `ScriptKittyMessageSchema`.
    - **Technical:** Low-latency network I/O, secure API access to external platforms, high concurrency to handle simultaneous user inputs.
    - **Data:** Raw text input from user, platform-specific message metadata.
  - **Dependencies (Dx):**
    - **Internal (DIx):**
      - **Script Kitty.Foundry (Kubernetes):** For deployment as horizontally scalable microservices, load balancing, and container runtime.
      - **Script Kitty.Foundry (Security):** For managing API keys to external platforms (via HashiCorp Vault).
    - **External (DEx):**
      - **External Messaging Platform APIs/SDKs:** Slack Events API, Telegram Bot API, custom WebSocket servers for web chat, OpenAPI/REST API endpoints.
  - **Correlations (Cx):**
    - **To Nexus's Input Normalization & Security Filter:** Directly feeds its output (`ScriptKittyMessageSchema`) to this component, ensuring normalized, initial security checks.

- **To Foundry's Observability:** Emits metrics (e.g., `messages_received_total`), logs (e.g., `incoming_message_payload`), and traces (for initial request span) that Foundry consumes.
- **To User Load:** Scales dynamically (HPA via Foundry) based on incoming message volume.

- **1.2. Input Normalization & Security Filter:**

  - **Requirements (Rx):**
    - **Functional:** Unicode normalization, whitespace cleanup, basic spell correction, preliminary PII redaction, heuristic toxicity/spam detection.
    - **Technical:** CPU for text processing, fast string manipulation capabilities.
    - **Data:** `ScriptKittyMessageSchema` object (specifically the `raw_text` field).
  - **Dependencies (Dx):**
    - **Internal (DIx):**
      - **Script Kitty.Foundry (Kubernetes):** For microservice deployment.
    - **External (DEx):**
      - **Python Libraries:** `unidecode`, `ftfy`, lightweight spell-checking libraries, regex engines.
      - **Heuristic Models/Rule Sets:** Pre-defined rules for spam/toxicity detection.
  - **Correlations (Cx):**
    - **To Nexus's Session Manager:** Passes the cleaned `ScriptKittyMessageSchema` to update session history.
    - **To Nexus's Intent & Entity Extraction:** Provides a cleaner, more standardized input, improving NLU accuracy.
    - **To Script Kitty.Guardian (Policy Enforcement):** Potentially flags highly toxic inputs for Guardian's more rigorous policy engine, though this is a lightweight initial pass.
    - **To Foundry's Observability:** Logs filtered/redacted content and any detected spam/toxicity events.

**2. Dialogue State & Context Management Sub-system:**

- **2.1. Session Manager:**

  - **Requirements (Rx):**
    - **Functional:** Persistent storage of `session_context` (conversation history, entities, preferences, task state). High read/write throughput for active sessions. Reliable long-term storage.

- **Technical:** Low-latency key-value access (for active sessions), ACID compliance for long-term storage, high availability, data consistency.
- **Data:** `session_id` as key, structured JSON object as value (containing `conversation_history`, `extracted_entities`, `user_preferences`, `current_task_state`, `last_activity_timestamp`).
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Foundry (Kubernetes):** For deploying Session Manager microservice.
    - **Script Kitty.Foundry (Data Management):** Redis cluster for caching, PostgreSQL for persistent storage.
  - **External (DEx):**
    - **Redis:** For in-memory caching and active session management.
    - **PostgreSQL:** For long-term session archival and warm storage.
    - **ORMs/Database Drivers:** (e.g., SQLAlchemy) for Python-PostgreSQL interaction.
- **Correlations (Cx):**
  - **To Nexus's Contextual Memory & Retrieval:** Directly provides the `session_context` object.
  - **To Nexus's Intent & Entity Extraction:** The Session Manager updates the session context with newly extracted entities and refined task states.
  - **To Script Kitty.Core (Goal Decomposition):** Core will query/update task-related `current_task_state` within the session via Nexus's proxy for multi-turn interactions.
  - **To Foundry's Observability:** Emits metrics (e.g., `active_sessions_count`, `session_load_latency`), logs (session updates), and traces.
  - **To Foundry's Data Management:** Depends on Foundry for the provisioning and management of Redis and PostgreSQL instances.
- **2.2. Contextual Memory & Retrieval:**

  - **Requirements (Rx):**
    - **Functional:** Retrieve relevant conversation history (sliding window), current entities, and user preferences for LLM prompting. Update session context with new information.
    - **Technical:** Efficient in-memory access to `session_context` object.
    - **Data:** `session_context` object.
  - **Dependencies (Dx):**
    - **Internal (DIx):**
      - **Nexus's Session Manager:** Directly queries for `session_context`.
  - **Correlations (Cx):**

- **To Nexus's Intent & Entity Extraction (Nexus LLM):** Crucially injects the contextual information into the LLM prompt, enabling coherent, multi-turn understanding.
- **To Nexus's Response Generation (Nexus LLM):** Provides context for generating persona-consistent and contextually relevant responses.
- **To Nexus's Session Manager:** Initiates updates to the `session_context` after LLM processing (e.g., updating extracted entities).

**3. Core NLU & Intent Routing Sub-system:**

- **3.1. Intent & Entity Extraction (Nexus LLM):**

  - **Requirements (Rx):**
    - **Functional:** Accurate identification of user intent and extraction of key entities from natural language. Output structured JSON.
    - **Technical:** High-performance GPU/CPU inference for LLM, low latency, sufficient memory for model weights, robust inference server.
    - **Data:** Cleaned user input (`ScriptKittyMessageSchema`), contextual memory (`session_context`), system prompts, few-shot examples.
  - **Dependencies (Dx):**
    - **Internal (DIx):**
      - **Nexus's Input Normalization & Security Filter:** Provides the clean input.
      - **Nexus's Contextual Memory & Retrieval:** Provides conversational context.
      - **Script Kitty.Foundry (Kubernetes):** For deploying and scaling the LLM inference service (vLLM/Ollama).
      - **Script Kitty.Foundry (Data Management):** For storing/retrieving LLM model artifacts (e.g., `gguf` files for Ollama).
      - **Script Kitty.Foundry (CI/CD):** For deploying new versions of the LLM and its inference service.
    - **External (DEx):**
      - **LLM Model:** Quantized Mistral-7B/Llama 3-8B.
      - **LLM Inference Server:** vLLM or Ollama.
      - **Prompt Engineering Techniques:** For structuring input to the LLM.
  - **Correlations (Cx):**
    - **To Nexus's Semantic Router:** Outputs the identified intent and extracted entities, which the router uses for decision-making.
    - **To Nexus's Response Generation:** The classified intent and entities inform the response synthesis.
    - **To Script Kitty.Guardian (Evaluation & Training):** Anonymized LLM inputs and outputs (along with human feedback) are consumed by

Guardian for continuous fine-tuning of this LLM via Foundry's ML pipelines.
- **To Foundry's Observability:** Emits LLM inference latency, token usage, and error metrics.

- **3.2. Semantic Router & Fallback Mechanism:**

  - **Requirements (Rx):**
    - **Functional:** Intelligent routing of queries based on intent and confidence. Ability to distinguish between general chat, complex task requests, and information queries. Fallback mechanisms.
    - **Technical:** Fast decision logic, potentially a small, fast traditional classifier.
    - **Data:** Identified intent and entities from Nexus LLM, confidence scores.
  - **Dependencies (Dx):**
    - **Internal (DIx):**
      - **Nexus's Intent & Entity Extraction:** Primary source of intent/entity data.
      - **Script Kitty.Core (Internal Communication Hub):** If intent is `TASK_REQUEST` or `COMPLEX_INFORMATION_QUERY`, it routes the request to Core.
      - **Script Kitty.Guardian (Policy Enforcement):** For very high-risk intents detected, might trigger a direct review/block by Guardian.
    - **External (DEx):**
      - **Hugging Face Transformers (Optional):** For a pre-trained classifier (e.g., `bert-base-uncased`) for fast-path intents.
      - **Custom Rules Engine:** For hardcoded routing rules.
  - **Correlations (Cx):**
    - **To Script Kitty.Core (Goal Decomposition):** *Crucial correlation.* This is the primary gateway for user requests to enter Core's domain, directly impacting Core's planning engine. A well-classified intent here reduces Core's initial ambiguity.
    - **To Nexus's Response Synthesis:** If Nexus handles the query directly (e.g., `GENERAL_CHAT`), it signals the response generation layer.
    - **To Foundry's Observability:** Logs routing decisions, fallback activations, and potential misclassifications.

## 4. Response Synthesis & Persona Layer:

- **4.1. Response Generation (Nexus LLM):**

  - **Requirements (Rx):**
    - **Functional:** Generate natural, coherent, and persona-consistent (Script Kitty's expert persona) chat responses. Incorporate structured `response_data` into natural language.

- **Technical:** High-performance GPU/CPU inference, low latency, robust inference server.
- **Data:** Original user query, identified intent/entities, `session_context`, `response_data` from Core/agents, system prompts for persona.
  - **Dependencies (Dx):**
    - **Internal (DIx):**
      - **Nexus's Intent & Entity Extraction:** Provides intent/entities.
      - **Nexus's Contextual Memory & Retrieval:** Provides conversational context.
      - **Nexus's Intermediary & Internal Communication Proxy:** Receives `response_data` from Core.
      - **Script Kitty.Foundry (Kubernetes):** For deploying and scaling the LLM inference service (vLLM/Ollama).
      - **Script Kitty.Foundry (Data Management):** For storing/retrieving LLM model artifacts.
    - **External (DEx):**
      - **LLM Model:** Same as NLU (Mistral-7B/Llama 3-8B).
      - **LLM Inference Server:** vLLM or Ollama.
      - **Prompt Engineering Techniques:** For persona injection and structured data integration.
  - **Correlations (Cx):**
    - **To Nexus's Group Chat Facilitation & Multi-Agent Attribution:** Generates the core text that this component then attributes and formats.
    - **To User Experience:** Directly impacts the perceived intelligence and helpfulness of Script Kitty.
    - **To Script Kitty.Guardian (Evaluation & Training):** Generated responses (and human feedback on them) are crucial for evaluating and fine-tuning this LLM.
- **4.2. Group Chat Facilitation & Multi-Agent Attribution:**

  - **Requirements (Rx):**
    - **Functional:** Prefix responses with agent attributions (e.g., "Armory reports:"). Manage turn-taking and progress updates from backend agents.
    - **Technical:** String manipulation, simple logic.
    - **Data:** Generated `response_text`, `source_agent` metadata, `event_type` (e.g., `PROGRESS_UPDATE`).
  - **Dependencies (Dx):**
    - **Internal (DIx):**
      - **Nexus's Response Generation:** Provides the raw generated text.
      - **Nexus's Intermediary & Internal Communication Proxy:** Receives `source_agent` and `event_type` metadata from Core.
    - **External (DEx):**

- **Templating Engines (e.g., Jinja2):** For dynamic insertion of data into predefined conversational templates.
  - **Correlations (Cx):**
    - **To User Experience:** This component is critical for maintaining the "group chat" illusion and transparency about which internal agent is performing which action.
    - **To Script Kitty.Core (Agent Orchestration):** Core's decision on which agent handles a sub-task directly influences the attribution prefix used here.

## 5. Intermediary & Internal Communication Proxy:

- **5.1. Core Communication Gateway (gRPC Client):**

  - **Requirements (Rx):**
    - **Functional:** High-performance, reliable, bidirectional communication with Script Kitty.Core. Efficient serialization/deserialization.
    - **Technical:** gRPC client libraries, Protobuf definition enforcement, robust connection management (retries, pooling).
    - **Data:** `ScriptKittyMessageSchema` (to Core), `TaskUpdateSchema`, `TaskResultSchema` (from Core).
  - **Dependencies (Dx):**
    - **Internal (DIx):**
      - **Nexus's Semantic Router:** Initiates outbound requests to Core.
      - **Nexus's Response Synthesis:** Receives inbound responses from Core.
      - **Script Kitty.Core (Internal Communication Hub):** The gRPC server endpoint for Core.
      - **Script Kitty.Foundry (Kubernetes/Service Mesh):** For network connectivity and service discovery to Core.
    - **External (DEx):**
      - **gRPC libraries:** For the chosen programming language (e.g., `grpcio` for Python).
      - **Protobuf compilers:** For generating language-specific message classes from `.proto` definitions.
  - **Correlations (Cx):**
    - **To Script Kitty.Core (Internal Communication Hub):** This is the direct communication channel. Any latency or failure here directly impacts Core's ability to receive user requests and Nexus's ability to get results.
    - **To Foundry's Observability:** Emits gRPC call latency, connection errors, and message counts.
- **5.2. User Output Renderer & Channel Dispatcher:**

  - **Requirements (Rx):**

- **Functional:** Convert internal `ScriptKittyOutputSchema` to platform-specific message format. Dispatch message to the correct channel.
- **Technical:** Low-latency API calls to external platforms, adherence to platform rate limits.
- **Data:** `ScriptKittyOutputSchema` (containing `response_text`, `suggested_actions`, `channel_type`).
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Nexus's Response Synthesis:** Provides the final output.
    - **Nexus's Channel Ingestion Adapters (in reverse):** Reuses these adapters for dispatching messages.
    - **Script Kitty.Foundry (Kubernetes):** For deploying dispatching logic.
  - **External (DEx):**
    - **External Messaging Platform APIs/SDKs:** Same as Channel Ingestion (Slack, Telegram, Web Sockets).
- **Correlations (Cx):**
  - **To User Experience:** The final touchpoint for the user, directly impacting the responsiveness and quality of the interaction.
  - **To Foundry's Observability:** Logs outbound message delivery status, API call errors, and rate limit hits.

---

## II. Script Kitty.Core - The Central Governing AI / High-Level Orchestrator

**Core Purpose (Recap):** The strategic brain. Receives high-level goals from Nexus, breaks them into executable plans, orchestrates agents, manages global state, integrates feedback.

**1. Goal Decomposition & HTN Planning Engine:**

- **Requirements (Rx):**
  - **Functional:** Interpret natural language goals, generate hierarchical task networks (HTN) or action sequences. Leverage global knowledge, adapt plans dynamically.
  - **Technical:** Powerful LLM inference (potentially larger), logical reasoning capabilities, access to knowledge graph, potentially symbolic planning engine.
  - **Data:** `TaskRequestSchema` from Nexus, `GlobalKnowledgeGraph` (GKGS), plan templates, historical execution data from Guardian.
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Nexus (Core Communication Gateway):** Receives incoming `TaskRequestSchema`.

- **Script Kitty.Core (Global Knowledge Graph Store - GKGS):** Reads knowledge, rules, and plan templates from GKGS. Writes refined plans and learned heuristics to GKGS.
- **Script Kitty.Core (Agent Orchestration):** Provides the decomposed plan for execution.
- **Script Kitty.Guardian (Feedback & Learning Integration):** Consumes plan success/failure metrics, feeds back into planning heuristics.
- **Script Kitty.Armory (Dynamic Tool Registry):** Queries Armory for available tools and their capabilities during planning.
- **Script Kitty.Skills (Specialized AI Models):** Queries Skills for available model capabilities during planning.
- **Script Kitty.Foundry (Kubernetes/LLM Inference):** For deploying and scaling the Core LLM.
  - **External (DEx):**
    - **Larger LLM (e.g., Llama 3-70B, GPT-4o API):** Via API adapter or self-hosted (managed by Foundry's LLM serving).
    - **Planning Libraries:** LangChain/AutoGen (for defining agent workflows), custom HTN planner implementations.
- **Correlations (Cx):**
  - **From Nexus:** The quality of intent classification and entity extraction by Nexus directly impacts the ambiguity Core's LLM needs to resolve, affecting planning latency and success.
  - **To Agent Orchestration:** *Crucial correlation.* Directly feeds the generated plan, dictating which agents are summoned and what tasks they are given. A robust plan leads to efficient agent execution.
  - **To Global Knowledge Graph:** Updates the GKGS with new planning strategies, learned task breakdowns, and improved templates based on successful executions.
  - **From Guardian:** Guardian's feedback on task success/failure directly refines Core's planning heuristics and plan templates, driving self-improvement.
  - **To Foundry's Observability:** Logs planning steps, LLM calls, and planning duration.

## 2. Agent Orchestration & Lifecycle Management:

- **Requirements (Rx):**
  - **Functional:** Dynamically deploy/activate, allocate tasks to, monitor progress of, and terminate specialized agents (Armory, Skills, Guardian instances).
  - **Technical:** Kubernetes API client access, workflow management, state tracking of active agents.
  - **Data:** Decomposed plan, agent registry (from GKGS), task updates from agents.
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Core (Goal Decomposition):** Receives the executable plan.

- - - **Script Kitty.Core (Internal Communication Hub):** For sending tasks to and receiving updates from agents.
    - **Script Kitty.Core (Global Knowledge Graph Store):** Reads agent registry (mapping agent types to Docker images/resource requirements).
    - **Script Kitty.Foundry (Kubernetes API Client):** Directly interacts with Kubernetes to deploy, scale, and manage agent pods.
  - **External (DEx):**
    - **Kubernetes Python Client:** For programmatic interaction with the Kubernetes API.
    - **Argo Workflows / Kubeflow Pipelines:** For defining and executing complex multi-agent workflows.
- **Correlations (Cx):**
  - **To Armory, Skills, Guardian:** *Crucial correlation.* Directly controls the instantiation and task allocation to these agents, dictating their workload and execution.
  - **From Armory, Skills, Guardian:** Receives `TaskUpdate` messages (progress, status, results), which it uses to update the plan's execution state.
  - **To Nexus (Core Communication Gateway):** Sends `TaskUpdate` messages back to Nexus for user display (via Nexus's Response Synthesis).
  - **To Foundry's Observability:** Logs agent lifecycle events (spawn, terminate), task allocation, and progress.

## 3. Global Knowledge Graph Store (GKGS) & Context Store:

- **Requirements (Rx):**
  - **Functional:** Central, authoritative, persistent storage for all Script Kitty's knowledge (facts, rules, learned heuristics, tool manifests, task contexts). Semantic search capabilities.
  - **Technical:** Graph database with strong query language, vector indexing for semantic search, high availability, data consistency.
  - **Data:** Nodes (entities), Edges (relationships), Properties (attributes). Plan templates, agent capabilities, success heuristics, safety rules, environmental facts, entity embeddings.
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Core (Goal Decomposition):** Reads planning knowledge, writes learned heuristics.
    - **Script Kitty.Core (Agent Orchestration):** Reads agent registry.
    - **Script Kitty.Armory (Dynamic Tool Registry):** Armory populates tool manifests into the GKGS.
    - **Script Kitty.Skills (Specialized AI Models):** Skills provides metadata on its models/capabilities for GKGS.
    - **Script Kitty.Guardian (Feedback & Learning Integration):** Reads/writes performance data, evaluation results, ethical policies.

- **Script Kitty.Foundry (Kubernetes/Data Management):** For deploying and managing the graph database and vector store.
  - **External (DEx):**
    - **Graph Database:** Neo4j Community Edition or ArangoDB.
    - **Vector Database:** Weaviate or Chroma.
    - **Knowledge Graph Embedding Libraries:** `pykeen`, `OpenKE`.
    - **RAG Libraries:** LlamaIndex, Haystack.
- **Correlations (Cx):**
  - **To Core (Goal Decomposition):** *Crucial correlation.* The GKGS directly informs Core's planning decisions by providing context, rules, and available tools. Any update to the GKGS (e.g., a new tool registered by Armory) immediately enriches Core's capabilities.
  - **From Armory (Dynamic Tool Registry):** Armory directly updates the GKGS with new tool manifests and capabilities.
  - **From Guardian (Automated Evaluation & Training Trigger):** Guardian updates the GKGS with evaluation results, learned policy refinements, and new training data references.
  - **To All Agents:** Conceptually, all agents can query the GKGS for global context, though often proxied through Core.
  - **To Foundry's Data Management:** Depends on Foundry for the persistent storage and scalable infrastructure of the graph database and vector store.

## 4. Internal Communication Hub & Message Broker:

- **Requirements (Rx):**
  - **Functional:** Manage structured message flow between Core and all backend agents. High-performance RPC for request/response. Asynchronous messaging for events.
  - **Technical:** gRPC server, Kafka/Pulsar cluster, Protobuf message schemas. High throughput, low latency, message durability.
  - **Data:** Standardized `Protobuf` messages (`ScriptKittyMessageSchema`, `TaskRequestSchema`, `TaskUpdateSchema`, `TaskResultSchema`, `AgentStatusSchema`).
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Nexus (Core Communication Gateway):** Receives from and sends to Nexus.
    - **Script Kitty.Core (Agent Orchestration):** Sends task requests to agents, receives updates.
    - **Script Kitty.Armory, Script Kitty.Skills, Script Kitty.Guardian:** These agents expose gRPC endpoints for Core to call, and push events to Kafka/Pulsar.
    - **Script Kitty.Foundry (Kubernetes/Service Mesh):** For network connectivity, service discovery, and message queue deployment.

- - ○ **External (DEx):**
      - **gRPC:** For synchronous RPC.
      - **Apache Kafka / Apache Pulsar:** For asynchronous event streaming.
      - **Protocol Buffers (Protobuf):** For all message schema definitions.
- **Correlations (Cx):**
  - **To All Agents:** *Crucial correlation.* This is the central nervous system's communication bus. Any disruption here halts Core's ability to orchestrate and receive updates.
  - **From All Agents:** All agents rely on this hub to receive tasks from Core and to send back progress updates, results, and events.
  - **To Foundry's Observability:** Monitors message queue depths, message processing latency, and gRPC call metrics across agents.

## 5. Feedback & Learning Integration:

- **Requirements (Rx):**
  - **Functional:** Ingest performance metrics, success/failure signals, human corrections, new observations. Refine planning heuristics, update GKGS, trigger system improvements.
  - **Technical:** Data pipeline for streaming data, MLflow for experiment tracking, DVC for data versioning, potentially RL components.
  - **Data:** `EvaluationResultSchema` from Guardian, `TaskSuccessMetrics`, `HumanFeedbackSchema`.
- **\*\*Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Guardian (Automated Evaluation & Training Trigger):** Guardian sends `EvaluationResultSchema` to Core.
    - **Script Kitty.Guardian (Human Oversight & Feedback Interface):** Provides human corrections.
    - **Script Kitty.Core (Global Knowledge Graph Store):** Writes refined heuristics and knowledge.
    - **Script Kitty.Foundry (Data Management):** For managing the data lake (for raw feedback data), DVC (for versioning datasets), and MLflow (for tracking learning experiments).
    - **Script Kitty.Foundry (ML Pipelines):** Core triggers pipelines managed by Foundry for retraining.
  - **External (DEx):**
    - **Apache Flink / Apache Spark Streaming:** For real-time processing of feedback streams.
    - **MLflow / DVC:** For tracking and versioning learning experiments.
    - **Reinforcement Learning Libraries:** (e.g., Ray RLlib - optional for policy optimization).
- **Correlations (Cx):**

- **From Guardian:** *Crucial correlation.* Guardian's output directly fuels Core's self-improvement loop. The quality and timeliness of Guardian's feedback dictate Core's learning speed.
  - **To Core (Goal Decomposition):** Directly refines the planning logic and templates within the Goal Decomposition engine. This is where Core's "intelligence" evolves.
  - **To Global Knowledge Graph:** Updates the GKGS with new and improved planning policies and heuristics.
  - **To Foundry (ML Pipelines):** Core initiates training pipelines in Foundry based on feedback, consuming data versioned by DVC and tracked by MLflow.

---

# III. Script Kitty.Armory - Resource Acquisition & Tooling AI

**Core Purpose (Recap):** Acquiring, processing external information, managing dynamic tool registry, and providing access to external resources.

## 1. Intelligent Web Scraping & Data Digestion:

- **Requirements (Rx):**
  - **Functional:** Autonomously navigate web, extract structured/unstructured data, summarize, perform NER/RE. Process into usable formats.
  - **Technical:** Headless browser automation, robust HTML parsing, NLP capabilities, proxy management (to avoid IP blocking), significant network I/O.
  - **Data:** Research queries (from Core), raw HTML/text, extracted entities, summaries.
- **Dependencies (Dx):
  - **Internal (DIx):**
    - **Script Kitty.Core (Internal Communication Hub):** Receives research queries from Core.
    - **Script Kitty.Core (Global Knowledge Graph Store):** Potentially reads existing knowledge before performing redundant research. Writes extracted entities/facts back to GKGS.
    - **Script Kitty.Foundry (Kubernetes):** For deploying web scraping agents as scalable pods.
    - **Script Kitty.Foundry (Security):** For managing external proxy credentials (via Vault).
    - **Script Kitty.Foundry (Data Management):** For storing raw and processed scraped data in the data lake.
  - **External (DEx):**
    - **Web Resources:** Any public/private websites, APIs.
    - **Headless Browsers/Automation:** Playwright, Selenium.
    - **HTML Parsers:** BeautifulSoup4, lxml.
    - **Text Extraction:** Trafilatura, boilerpy3.

- - **NLP Libraries:** spaCy, Hugging Face Transformers (for summarization, NER).
    - **Proxies/VPNs:** To bypass geo-restrictions or IP blocks.
  - **Correlations (Cx):**
    - **From Core (Goal Decomposition):** Core issues specific research tasks, which directly drive Armory's web scraping activities.
    - **To Core (Global Knowledge Graph Store):** Armory enriches the GKGS with newly acquired and processed information, directly expanding Script Kitty's knowledge base and thus Core's planning capabilities.
    - **To Skills (Data Processing & Analysis):** Processed data might be handed off to Skills for deeper analysis.
    - **To Foundry's Observability:** Logs successful/failed scrapes, data volume, and network usage.

## 2. Dynamic Tool Registry & Discovery:

- **Requirements (Rx):**
  - **Functional:** Maintain searchable registry of all internal/external tools. Provide detailed manifests (capabilities, params, schemas). Discover tools based on semantic capabilities.
  - **Technical:** Database for metadata, search engine for semantic search, efficient tool onboarding.
  - **Data:** `ToolManifestSchema` (tool name, description, capabilities, inputs, outputs, adapter type, environment requirements).
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Core (Global Knowledge Graph Store):** *Crucial correlation.* This sub-system *populates* the GKGS with tool manifests. Core (specifically its planning engine) *reads* from the GKGS to discover available tools.
    - **Script Kitty.Foundry (Data Management):** PostgreSQL for structured storage, Elasticsearch/OpenSearch for search.
    - **Script Kitty.Foundry (CI/CD):** For automating ingestion of new tool definitions.
  - **External (DEx):**
    - **PostgreSQL:** For structured storage of tool manifests.
    - **Elasticsearch/OpenSearch:** For semantic search over tool capabilities.
    - **OpenAPI/Swagger:** For defining standardized REST API schemas for external tools.
    - **Custom Parsers:** For ingesting tool documentation (man pages, source code comments).
- **Correlations (Cx):**

- **To Core (Goal Decomposition):** Directly enables Core to identify and integrate relevant tools into its plans. A newly registered tool in Armory immediately becomes available to Core's planning.
- **To Skills (Sandboxed Code Execution):** If a tool requires execution, Armory's manifest tells Core how to invoke it, which then passes to Skills for execution.
- **From Armory (Templated API/CLI Wrapper Generation):** New wrappers created here are registered in this registry.

## 3. Templated API/CLI Wrapper Generation (Human-Verified):

- **Requirements (Rx):**
  - **Functional:** Generate secure, templated code wrappers for novel or custom tools. Present for human review.
  - **Technical:** Code-focused LLM inference, secure templating, static code analysis, web UI for human review, version control integration.
  - **Data:** Tool API/CLI documentation, wrapper templates, LLM-generated code, human review feedback.
- **Dependencies (Dx):
  - **Internal (DIx):**
    - **Script Kitty.Core (Internal Communication Hub):** Receives requests for wrapper generation from Core (if Core identifies a need).
    - **Script Kitty.Armory (Dynamic Tool Registry):** Registers newly approved wrappers.
    - **Script Kitty.Foundry (Kubernetes):** For deploying LLM inference for code generation.
    - **Script Kitty.Foundry (CI/CD):** For running static code analysis and integrating with Git for human review.
    - **Script Kitty.Foundry (Security):** For providing secure sandboxed environment for testing generated code if needed, and managing access to human review UI.
  - **External (DEx):**
    - **Code LLM:** CodeLlama, StarCoder2, DeepSeek Coder (served by vLLM/TGI).
    - **Templating Engines:** Jinja2.
    - **Static Code Analysis Tools:** Bandit (Python), ESLint (JavaScript), SonarQube Community Edition.
    - **Git/GitLab/GitHub:** For version control and pull request-based human review workflow.
    - **Web UI Frameworks:** React/Next.js (frontend), FastAPI (backend).
- **Correlations (Cx):**
  - **To Armory (Dynamic Tool Registry):** A successfully generated and human-approved wrapper is immediately added to the tool registry, expanding Armory's operational capabilities.

- **From Core (Goal Decomposition):** Core might identify a need for a new tool wrapper (e.g., if it attempts to use a capability for which no wrapper exists), triggering this process.
- **To Script Kitty.Guardian (Human Oversight):** Requires human review and approval, making Guardian's human-in-the-loop interface a critical dependency for security and correctness.

**4. Resource Provisioning via IaC Templates:**

- **Requirements (Rx):**
  - **Functional:** Programmatically provision cloud VMs, GPUs, quantum computing access, storage. Utilize pre-approved IaC templates.
  - **Technical:** Access to cloud provider APIs, IaC tool integration, secure credential management.
  - **Data:** IaC blueprints (Terraform HCL, Ansible Playbooks), provisioning parameters (from Core).
- **\*\*Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Core (Agent Orchestration):** Receives requests from Core to provision specific resources for tasks (e.g., a high-GPU instance for a Skills agent, a quantum compute job).
    - **Script Kitty.Foundry (Security):** For managing cloud API keys (via Vault).
    - **Script Kitty.Foundry (Kubernetes API Client):** For requesting ephemeral Kubernetes pods with custom resource requests.
  - **External (DEx):**
    - **Terraform:** For declarative infrastructure provisioning.
    - **Ansible:** For configuration management.
    - **Cloud Provider SDKs:** (e.g., `boto3` for AWS, `google-cloud-python` for GCP).
    - **Quantum Computing SDKs:** Qiskit, Cirq, PennyLane.
- **Correlations (Cx):**
  - **To Core (Agent Orchestration):** Core relies on Armory to provision the necessary compute environments for Skills agents to execute their tasks.
  - **To Skills (Specialized AI Models & Sandboxed Code Execution):** Provides the underlying hardware/software environment for Skills to operate.
  - **To Foundry's Observability:** Logs resource provisioning events, success/failure, and cost metrics.
  - **To Script Kitty.Guardian (Policy Enforcement):** Resource requests might be subject to Guardian's policies (e.g., budget limits, allowed regions).

---

# IV. Script Kitty.Skills - Task Execution AI

**Core Purpose (Recap):** Houses specialized AI models and agents that execute specific computational tasks identified by Script Kitty.Core. These are the "doers."

## 1. Specialized AI Models & Inference:

- **Requirements (Rx):**
  - **Functional:** Run diverse pre-trained/fine-tuned ML models (CV, NLP, time-series, code generation). Perform inference on specific data.
  - **Technical:** Optimized model serving frameworks, efficient GPU/CPU utilization, access to model artifacts.
  - **Data:** Input data for inference (from Core), model weights, inference results.
- **\*\*Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Core (Internal Communication Hub):** Receives task requests from Core (specifying model to use, input data). Sends inference results back to Core.
    - **Script Kitty.Armory (Resource Provisioning):** May require Armory to provision specific GPU resources for large models.
    - **Script Kitty.Foundry (Kubernetes):** For deploying model serving microservices as scalable pods.
    - **Script Kitty.Foundry (Data Management):** For retrieving model artifacts (from DVC-versioned storage/feature store).
    - **Script Kitty.Foundry (CI/CD):** For deploying new/retrained model versions.
  - **External (DEx):**
    - **ML Frameworks:** PyTorch, TensorFlow, Keras, scikit-learn.
    - **Model Serving Frameworks:** TorchServe, TensorFlow Serving, KServe, Ray Serve, ONNX Runtime.
    - **Hugging Face Transformers:** For access to pre-trained models.
    - **Specialized Models:** Domain-specific models (e.g., medical imaging, financial forecasting).
- **Correlations (Cx):**
  - **From Core (Agent Orchestration):** Core directly instructs Skills to run specific models with specific inputs.
  - **To Core (Internal Communication Hub):** Returns results of inference, allowing Core to integrate them into its plan.
  - **From Guardian (Automated Evaluation & Training):** Guardian sends updated/retrained models to Foundry, which are then deployed to Skills. Guardian also monitors Skills' model performance.
  - **To Foundry's Observability:** Emits inference latency, throughput, GPU utilization, and model confidence metrics.

## 2. Sandboxed Code Execution & CLI Interaction:

- **Requirements (Rx):**

- ○ **Functional:** Safely execute arbitrary code (Python, R, shell scripts) or CLI commands within an isolated environment. Strict resource limits.
  - ○ **Technical:** Containerization, Linux namespaces/cgroups, syscall filtering, network isolation. Fast spin-up/tear-down of sandboxes.
  - ○ **Data:** Code to execute, CLI commands, input files, execution output (stdout/stderr), exit codes.
- ● **Dependencies (Dx):
  - ○ **Internal (DIx):**
    - ■ **Script Kitty.Core (Internal Communication Hub):** Receives code/commands from Core. Returns execution results.
    - ■ **Script Kitty.Armory (Templated API/CLI Wrapper Generation):** Executes wrappers generated by Armory.
    - ■ **Script Kitty.Foundry (Kubernetes):** For creating ephemeral Docker containers/pods for sandboxing.
    - ■ **Script Kitty.Foundry (Security):** *Crucial dependency.* Relies heavily on Foundry's container security (Docker, Linux Namespaces/cgroups, seccomp-bpf, Firejail/Bubblewrap) and network policies (Calico/Cilium) to prevent escapes or malicious activity.
    - ■ **Script Kitty.Foundry (Data Management):** For secure temporary storage of input/output files.
  - ○ **External (DEx):**
    - ■ **Docker:** For containerization.
    - ■ **Linux Kernel Features:** Namespaces, cgroups, seccomp-bpf.
    - ■ **User-space Sandboxing:** Firejail, Bubblewrap.
    - ■ **Jupyter Kernel Gateway/Binder:** For interactive or complex script execution environments.
    - ■ `subprocess` **module (Python):** For running CLI commands.
- ● **Correlations (Cx):**
  - ○ **From Core (Agent Orchestration):** Core sends execution tasks to Skills, directly triggering code execution.
  - ○ **To Core (Internal Communication Hub):** Sends back execution results, logs, and any errors.
  - ○ **To Script Kitty.Guardian (Policy Enforcement):** Guardian monitors and logs all sandboxed execution attempts and results, flagging any policy violations (e.g., attempts to access restricted files, excessive resource consumption). This is a critical security feedback loop.
  - ○ **To Foundry's Observability:** Emits metrics on execution duration, resource consumption, and security alerts (from Falco integration).

## 3. Data Processing & Analysis Frameworks:

- ● **Requirements (Rx):**
  - ○ **Functional:** Perform data manipulation, cleaning, statistical analysis, and visualization. Handle large datasets.

- ○ **Technical:** High memory/CPU throughput, potentially distributed processing capabilities.
  - ○ **Data:** Raw/processed data from Armory/Core, analysis results, visualizations.
- ● **Dependencies (Dx):
  - ○ **Internal (DIx):**
    - ■ **Script Kitty.Core (Internal Communication Hub):** Receives data processing tasks from Core. Returns analysis results.
    - ■ **Script Kitty.Armory (Web Scraping & Data Digestion):** Consumes data from Armory for analysis.
    - ■ **Script Kitty.Foundry (Data Management):** Accesses data lake (for large datasets), feature store (for pre-computed features).
    - ■ **Script Kitty.Foundry (Kubernetes):** For deploying scalable data processing pods.
  - ○ **External (DEx):**
    - ■ **Python Libraries:** Pandas, NumPy, SciPy, Matplotlib, Seaborn, Plotly.
    - ■ **Distributed Processing:** Dask, Apache Spark (PySpark).
    - ■ **High-performance columnar:** Polars, DuckDB.
- ● **Correlations (Cx):**
  - ○ **From Core (Agent Orchestration):** Core's plan dictates what data analysis Skills performs.
  - ○ **To Core (Internal Communication Hub):** Provides structured analysis results that Core can use for decision-making or present to the user via Nexus.
  - ○ **To Script Kitty.Guardian (Evaluation):** Analysis results might be used by Guardian to evaluate outcomes or identify new learning opportunities.

## 4. Generative Task Execution (e.g., Programming, Report Generation):

- ● **Requirements (Rx):**
  - ○ **Functional:** Generate code, reports, or complex textual outputs based on instructions. Coherent, accurate generation.
  - ○ **Technical:** Code-focused LLM inference, potentially domain-specific knowledge integration, syntax/semantic validation.
  - ○ **Data:** Task instructions (from Core), context, generated code/text.
- ● **Dependencies (Dx):
  - ○ **Internal (DIx):**
    - ■ **Script Kitty.Core (Internal Communication Hub):** Receives generative tasks. Sends back generated content.
    - ■ **Script Kitty.Core (Global Knowledge Graph Store):** Potentially queries GKGS for relevant domain knowledge.
    - ■ **Script Kitty.Foundry (Kubernetes):** For deploying generative LLM inference services.
    - ■ **Script Kitty.Foundry (Data Management):** For storing/retrieving fine-tuned models.
  - ○ **External (DEx):**
    - ■ **Code-focused LLMs:** CodeLlama, StarCoder2, DeepSeek Coder.

- - **Text Generation LLMs:** For reports.
  - **LLM Inference Servers:** vLLM, Text Generation Inference.
  - **Code Parsers/Validators:** `ast` module (Python), `tree-sitter`.
  - **Unit Testing Frameworks:** `pytest` (for self-generated tests).
- **Correlations (Cx):**
  - **From Core (Agent Orchestration):** Core specifies the programming or generation task.
  - **To Core (Internal Communication Hub):** Delivers the generated code/report back to Core for review, further action (e.g., execution by Skills' sandboxed environment), or presentation to the user.
  - **To Script Kitty.Guardian (Evaluation):** Generated code/reports are subject to evaluation by Guardian (e.g., code correctness, report quality, safety review). Positive/negative feedback is used for retraining.
  - **To Foundry's Observability:** Emits metrics on generation latency, token usage, and potentially code quality metrics.

---

# V. Script Kitty.Guardian - Safety & Alignment Proxy / Evaluation & Training AI

**Core Purpose (Recap):** The system's conscience and continuous improvement engine. Enforces controls, monitors performance, identifies learning opportunities, triggers training.

## 1. Policy Enforcement Engine:

- **Requirements (Rx):**
  - **Functional:** Intercept proposed actions from Core/agents. Evaluate against ethical/legal/safety policies. Flag violations. Provide explanations.
  - **Technical:** Rule engine, LLM for ethical context, XAI tools, low latency for critical path decisions.
  - **Data:** Proposed `ActionRequestSchema` (from Core/agents), policy rules, ethical guidelines, human feedback.
- **\*\*Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Core (Internal Communication Hub):** Core sends proposed actions to Guardian for pre-execution policy checks. Guardian sends back approval/denial.
    - **Script Kitty.Armory, Script Kitty.Skills:** Their proposed actions are routed through Guardian's policy engine.
    - **Script Kitty.Core (Global Knowledge Graph Store):** Reads ethical policies, safe action patterns, and contextual safety information from GKGS. Writes policy refinements.

- ■ **Script Kitty.Foundry (Kubernetes):** For deploying policy engine as a critical, high-availability microservice.
- ■ **Script Kitty.Foundry (Security):** Policy engine itself relies on Foundry's security (network policies, access control).
  - ○ **External (DEx):**
    - ■ **Open Policy Agent (OPA):** For declarative policy definition (Rego language).
    - ■ **Custom Rule Engine:** For complex ethical dilemmas.
    - ■ **Ethical LLM (fine-tuned):** For identifying ethical conflicts.
    - ■ **Explainable AI (XAI) Tools:** LIME, SHAP (for explaining LLM/policy decisions).
- ● **Correlations (Cx):**
  - ○ **To Core & Other Agents:** *Crucial correlation.* Guardian's policy engine acts as a gatekeeper. Any proposed action must pass through it, directly influencing the behavior and capabilities of all other agents. A denial here can halt a task.
  - ○ **To Human Oversight & Feedback Interface:** Flags policy violations for human review and override, creating a critical human-in-the-loop mechanism.
  - ○ **To Foundry's Observability:** Logs all policy evaluations, violations, and human overrides, providing an audit trail.

**2. Real-time Performance Monitoring & Anomaly Detection:**

- ● **Requirements (Rx):**
  - ○ **Functional:** Continuously collect metrics from all Script Kitty components. Detect performance degradation, error rate spikes, abnormal resource usage, model confidence drops.
  - ○ **Technical:** Stream processing, anomaly detection algorithms, integration with Foundry's observability stack.
  - ○ **Data:** Metrics streams (from Prometheus), log streams (from Loki/Elasticsearch), trace data (from Jaeger/Zipkin).
- ● **\*\*Dependencies (Dx):**
  - ○ **Internal (DIx):**
    - ■ **Script Kitty.Foundry (Observability Stack):** *Crucial dependency.* Guardian is the primary consumer of the aggregated metrics, logs, and traces provided by Foundry. Without Foundry's robust observability, Guardian cannot perform its monitoring role.
  - ○ **External (DEx):**
    - ■ **Anomaly Detection Libraries:** PyOD, Isolation Forest, or custom ML models.
    - ■ **Stream Processing Frameworks:** Apache Flink, Apache Spark Streaming (for real-time analysis of metrics/logs).
- ● **Correlations (Cx):**
  - ○ **From All Components (via Foundry):** Receives a continuous stream of operational data, forming its primary input.

- ○ **To Automated Evaluation & Training Trigger:** Detected anomalies or performance degradation directly trigger evaluation and potential retraining pipelines.
- ○ **To Human Oversight & Feedback Interface:** Alerts human operators to critical system issues.

### 3. Automated Evaluation & Training Trigger:

- **Requirements (Rx):**
  - ○ **Functional:** Automatically trigger retraining pipelines for specific models/modules based on performance, feedback, or schedule. Evaluate effectiveness of previous updates.
  - ○ **Technical:** Workflow orchestration, experiment tracking, data versioning, hyperparameter optimization.
  - ○ **Data:** Monitoring data, anomaly alerts, human feedback, success/failure signals, evaluation metrics.
- **\*\*Dependencies (Dx):**
  - ○ **Internal (DIx):**
    - ■ **Script Kitty.Guardian (Performance Monitoring):** Receives anomaly alerts and performance metrics.
    - ■ **Script Kitty.Guardian (Human Oversight):** Receives human feedback.
    - ■ **Script Kitty.Core (Feedback & Learning Integration):** Sends evaluation results to Core to refine planning.
    - ■ **Script Kitty.Foundry (MLOps & Infrastructure Fabric):** *Crucial dependency.* Relies heavily on Foundry's Kubeflow Pipelines/Argo Workflows for pipeline orchestration, MLflow for experiment tracking, DVC for data versioning, and Ray Tune/Optuna for hyperparameter optimization. Foundry *executes* the training, Guardian *triggers* and *evaluates* it.
    - ■ **Script Kitty.Foundry (Data Management):** Accesses the data lake for training datasets.
  - ○ **External (DEx):**
    - ■ **Kubeflow Pipelines / Argo Workflows:** For defining and executing ML pipelines.
    - ■ **MLflow:** For tracking experiments.
    - ■ **DVC:** For managing data versions.
    - ■ **Ray Tune / Optuna:** For hyperparameter optimization.
- **Correlations (Cx):**
  - ○ **To Foundry (ML Pipelines):** *Crucial correlation.* Directly initiates model retraining and deployment pipelines, making Guardian the driver of Script Kitty's autonomous learning.
  - ○ **To Core (Feedback & Learning Integration):** Provides evaluation results that inform Core's meta-learning and planning heuristic adjustments.
  - ○ **To Nexus, Core, Armory, Skills:** New model versions deployed through Foundry will directly impact the performance and behavior of these components.

**4. Bias Detection & Mitigation:**

- **Requirements (Rx):**
    - **Functional:** Analyze training data and model outputs for biases. Provide tools for human review. Suggest mitigation strategies.
    - **Technical:** Fairness metrics, bias detection algorithms, explainability tools.
    - **Data:** Training datasets, model outputs (especially from generative models from Nexus/Skills).
- **\*\*Dependencies (Dx):**
    - **Internal (DIx):**
        - **Script Kitty.Foundry (Data Management):** Accesses training datasets and model outputs from the data lake/object storage.
        - **Script Kitty.Foundry (ML Pipelines):** Integrated into retraining pipelines.
        - **Script Kitty.Nexus (Response Generation):** Its outputs are analyzed.
        - **Script Kitty.Skills (Generative Task Execution):** Its outputs are analyzed.
        - **Script Kitty.Guardian (Human Oversight):** Flags detected biases for human review.
    - **External (DEx):**
        - **Fairness Toolkits:** Aequitas, Fairlearn, Google's What-If Tool, IBM's AI Fairness 360 (AIF360).
- **Correlations (Cx):**
    - **To Automated Evaluation & Training Trigger:** If significant bias is detected, it triggers a retraining cycle with bias mitigation strategies.
    - **To Human Oversight & Feedback Interface:** Presents bias reports for human review and intervention.
    - **To All LLM/Generative Components (Nexus, Skills):** Findings from bias detection directly influence their fine-tuning parameters and data preparation, leading to more ethical outputs.

**5. Human Oversight & Feedback Interface:**

- **Requirements (Rx):**
    - **Functional:** Provide clear summaries of policy violations, performance issues, ethical dilemmas. Allow human approval, override, corrective feedback.
    - **Technical:** Interactive web dashboard, secure authentication, API for feedback submission.
    - **Data:** Policy violation alerts, performance reports, bias reports, human approval/override/feedback.
- **\*\*Dependencies (Dx):**
    - **Internal (DIx):**
        - **Script Kitty.Guardian (Policy Enforcement):** Feeds policy violation alerts.
        - **Script Kitty.Guardian (Performance Monitoring):** Feeds system alerts.
        - **Script Kitty.Guardian (Bias Detection):** Feeds bias reports.

- - **Script Kitty.Core (Feedback & Learning Integration):** Receives human feedback.
    - **Script Kitty.Foundry (Security):** *Crucial dependency.* Keycloak for authentication/authorization for human users. Vault for managing UI secrets.
  - **External (DEx):**
    - **Web UI Frameworks:** Streamlit, Plotly Dash, React/Next.js (frontend), FastAPI (backend).
    - **Rasa Open Source:** For structured conversational feedback.
- **Correlations (Cx):**
  - **To Policy Enforcement Engine:** Human overrides directly influence the policy engine's behavior for specific cases, establishing a human-in-the-loop safety net.
  - **To Automated Evaluation & Training Trigger:** Human feedback is a direct input for triggering specific retraining cycles or prioritizing model improvements.
  - **To Core (Feedback & Learning Integration):** Human feedback on task success/failure or problematic outputs directly informs Core's planning refinement.
  - **To All Components:** Provides the ultimate human oversight and correction mechanism for the entire Script Kitty system's behavior and learning.

---

# VI. Script Kitty.Foundry - The MLOps & Infrastructure Fabric

**Core Purpose (Recap):** The foundational layer providing MLOps, CI/CD, data management, and compute orchestration. Underpins all other components.

**1. Kubernetes-Native Compute & Orchestration:**

- **Requirements (Rx):**
  - **Functional:** Provision, manage, scale, and heal all Script Kitty microservices and agents. Efficient resource allocation.
  - **Technical:** Robust Kubernetes cluster, autoscaling capabilities (HPA, VPA, Cluster Autoscaler), secure access to cloud provider APIs.
  - **Data:** Kubernetes manifests, resource requests/limits, pod definitions.
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Core (Agent Orchestration):** Core sends requests to Foundry to deploy/manage agents.
    - **All Other Script Kitty Components (Nexus, Core, Armory, Skills, Guardian):** *Crucial dependency.* Every single other component relies on Foundry for its existence and runtime environment. They are deployed as containers orchestrated by Foundry.
  - **External (DEx):**
    - **Kubernetes:** The core orchestration platform.

- **Cloud Provider APIs:** (e.g., AWS EC2, GCP Compute Engine) for provisioning underlying VMs.
- **KubeVirt:** (Optional) For VM management if needed.
- **Kubernetes Python Client:** For programmatic control.
- **Prometheus Operator/Grafana Operator:** For managing observability stack within Kubernetes.
- **Correlations (Cx):**
  - **To All Components:** *Fundamental correlation.* This is the lifeline. Its performance, stability, and scaling capabilities directly determine the performance, availability, and scalability of every other Script Kitty component.
  - **From Core (Agent Orchestration):** Core's orchestration decisions directly translate into deployment/scaling actions in Foundry.
  - **To Foundry's Observability:** Emits detailed metrics on cluster health, resource utilization, pod status, and autoscaling events.

**2. Container Registry & Image Management:**

- **Requirements (Rx):**
  - **Functional:** Securely store, version, and distribute Docker images for all Script Kitty components and models.
  - **Technical:** High availability, secure access, vulnerability scanning, content trust.
  - **Data:** Docker images, image tags, vulnerability reports, SBOMs.
- **Dependencies (Dx):**
  - **Internal (DIx):**
    - **Script Kitty.Foundry (CI/CD):** Pushes newly built images to the registry. Pulls images for deployment.
    - **All Other Script Kitty Components:** Their Docker images are stored here.
  - **External (DEx):**
    - **Harbor, Quay.io, Docker Hub, ECR, GCR:** Choice of container registry.
    - **Trivy, Clair:** For vulnerability scanning.
    - **Notary:** For image signing.
- **Correlations (Cx):**
  - **To Foundry's CI/CD:** Provides the central repository for deployable artifacts, directly feeding the CD pipeline.
  - **To Kubernetes (via Foundry):** Kubernetes pulls images from this registry to instantiate pods.
  - **To Foundry's Security:** Enforces image security policies.

**3. CI/CD & GitOps Automation:**

- **Requirements (Rx):**
  - **Functional:** Automate build, test, deployment, and update processes. Enforce Git as single source of truth. Automated rollbacks.

- ○ **Technical:** CI/CD pipeline orchestrator, Git integration, Kubernetes API access, declarative configuration.
  - ○ **Data:** Source code, test results, build artifacts, Kubernetes manifests, Git commit history.
- ● **Dependencies (Dx):**
  - ○ **Internal (DIx):**
    - ■ **All Script Kitty Components' Repositories:** Source code for Nexus, Core, Armory, Skills, Guardian, and Foundry itself.
    - ■ **Script Kitty.Foundry (Container Registry):** Pushes/pulls images.
    - ■ **Script Kitty.Foundry (Kubernetes):** Applies manifests to the cluster.
    - ■ **Script Kitty.Foundry (Observability):** Logs pipeline status, test results.
  - ○ **External (DEx):**
    - ■ **GitLab CI/CD, GitHub Actions, Jenkins, Tekton Pipelines:** CI/CD orchestration.
    - ■ **Argo CD, Flux CD:** GitOps agents.
    - ■ **Git:** Version control system (central for all code/configs).
    - ■ **Helm:** For packaging Kubernetes applications.
- ● **Correlations (Cx):**
  - ○ **To All Components:** *Crucial correlation.* This system dictates the speed, reliability, and safety of deploying new features, bug fixes, and model updates across the *entire* Script Kitty system. It directly impacts Script Kitty's "continuous evolution."
  - ○ **From Guardian (Automated Evaluation & Training Trigger):** Guardian triggers CD pipelines for model retraining and deployment.
  - ○ **To Kubernetes:** Argo CD/Flux CD continuously syncs the desired state (from Git) with the actual state in Kubernetes.

## 4. Observability Stack (Centralized Logging, Metrics, Tracing):

- ● **Requirements (Rx):**
  - ○ **Functional:** Collect, aggregate, store, and visualize logs, metrics, and traces from all services. Enable monitoring, debugging, performance analysis.
  - ○ **Technical:** High-throughput data ingestion, scalable storage, powerful querying, real-time visualization.
  - ○ **Data:** Logs, metrics, traces from all Script Kitty microservices and agents.
- ● **Dependencies (Dx):**
  - ○ **Internal (DIx):**
    - ■ **All Script Kitty Components:** *Crucial dependency for data emission.* Every component is instrumented to emit logs, metrics, and traces.
    - ■ **Script Kitty.Guardian (Performance Monitoring):** *Crucial dependency for consumption.* Guardian is the primary consumer of this data for anomaly detection and evaluation.
    - ■ **Script Kitty.Guardian (Human Oversight):** Provides dashboards for human operators.
  - ○ **External (DEx):**

- ■ **Fluent Bit/Fluentd:** Log collection agents.
- ■ **Loki, Elasticsearch/OpenSearch:** Log aggregation.
- ■ **Prometheus:** Metrics collection.
- ■ **Grafana:** Dashboards, alerting.
- ■ **OpenTelemetry SDKs:** For instrumentation across all services.
- ■ **Jaeger/Zipkin:** Trace visualization.
- ● **Correlations (Cx):**
  - ○ **From All Components:** Every component feeds data into the observability stack.
  - ○ **To Guardian:** *Crucial correlation.* The quality and completeness of observability data directly determines Guardian's ability to monitor, evaluate, and trigger learning. Without this data, Script Kitty is blind.
  - ○ **To Foundry (CI/CD):** Provides feedback on deployment health.
  - ○ **To SRE/DevOps Teams:** Enables manual debugging and performance tuning.

## 5. Data Management (Data Lake, Data Versioning, Feature Store):

- ● **Requirements (Rx):**
  - ○ **Functional:** Scalable storage for raw/processed data, model artifacts, features. Data versioning and reproducibility. Consistent feature serving.
  - ○ **Technical:** High-performance object storage, distributed file system capabilities, low-latency feature serving.
  - ○ **Data:** Raw inputs, scraped data, training datasets, model checkpoints, inference logs, features, feedback data.
- ● **Dependencies (Dx):**
  - ○ **Internal (DIx):**
    - ■ **Script Kitty.Armory (Web Scraping):** Stores raw and processed scraped data.
    - ■ **Script Kitty.Skills (Specialized AI Models):** Accesses training data and serves/consumes features.
    - ■ **Script Kitty.Guardian (Evaluation & Training):** Stores evaluation results, training datasets, model artifacts. Drives data versioning.
    - ■ **Script Kitty.Core (Global Knowledge Graph Store):** May store large knowledge assets.
  - ○ **External (DEx):**
    - ■ **MinIO, Ceph, AWS S3, GCP GCS:** Object storage.
    - ■ **DVC:** Data version control.
    - ■ **Feast:** Feature store (offline and online stores).
    - ■ **Apache Airflow, Prefect, Dagster:** Data orchestration.
- ● **Correlations (Cx):**
  - ○ **To Armory, Skills, Guardian:** Provides the fundamental data infrastructure for these components to operate. Data access is governed by this layer.
  - ○ **From Armory:** Feeds raw and processed external data into the data lake.
  - ○ **From Guardian:** Stores evaluation data and triggers retraining pipelines which consume data from here.

- ○ **To Foundry's ML Pipelines:** Provides versioned datasets and features for model training.

## 6. Security & Identity Management:

- ● **Requirements (Rx):**
  - ○ **Functional:** Centralized access control, secrets management, network policies, runtime security, container scanning. End-to-end security.
  - ○ **Technical:** Secure vault, policy enforcement agents, service mesh, identity provider, runtime monitoring.
  - ○ **Data:** Credentials, policies, audit logs, security alerts, user/role definitions.
- ● **Dependencies (Dx):**
  - ○ **Internal (DIx):**
    - ■ **All Script Kitty Components:** *Crucial dependency.* Every component relies on Foundry's security for safe operation (e.g., Vault for secrets, Calico for network isolation, Istio for mTLS).
    - ■ **Script Kitty.Guardian (Policy Enforcement):** Guardian's policy enforcement is often implemented using tools managed by Foundry (OPA).
    - ■ **Script Kitty.Guardian (Human Oversight):** Human login to dashboards is managed by Keycloak.
  - ○ **External (DEx):**
    - ■ **HashiCorp Vault:** Secrets management.
    - ■ **Calico, Cilium:** Network policy.
    - ■ **Istio, Linkerd:** Service mesh (mTLS, access control).
    - ■ **Falco:** Runtime security.
    - ■ **Trivy, Clair:** Container scanning.
    - ■ **Keycloak:** Identity & Access Management.
- ● **Correlations (Cx):**
  - ○ **To All Components:** *Fundamental correlation.* This layer directly dictates what every component can access, who can interact with it, and how securely it operates. A robust security posture here is non-negotiable for Script Kitty's integrity and trustworthiness.
  - ○ **From Guardian (Policy Enforcement):** Guardian defines high-level policies that Foundry's security tools (like OPA) then enforce.
  - ○ **To Foundry's Observability:** Emits detailed security audit logs and alerts, consumed by Guardian and SRE teams.

---

## Concluding Synthesis:

This exhaustive analysis of requirements, dependencies, and correlations reveals Script Kitty not as a collection of isolated modules, but as a deeply interconnected, symbiotic organism.

- **Foundry** is the **nervous system, skeletal structure, and circulatory system**, providing the essential infrastructure, communication pathways, security, and MLOps fabric that enables life. Every other component's very existence and operational capability are *direct requirements* met by Foundry.
- **Nexus** is the **sensory input and vocal output**, translating the external world for Core and presenting Core's actions back to the user. Its requirements for seamless interaction are dependent on Core's intelligence and Foundry's robust communication.
- **Core** is the **cerebral cortex**, responsible for strategic thought, planning, and orchestration. Its requirements for intelligent decision-making are critically dependent on the rich knowledge from GKGS, the tools registered by Armory, the execution capabilities of Skills, and the crucial feedback and guardrails provided by Guardian. Its outputs directly dictate the actions of Armory and Skills.
- **Armory** is the **limbs and hands for resource acquisition and tool integration**. Its requirements for external interaction and tool management are dependent on Core's directives and Foundry's secure external access. Its outputs (new data, new tools) directly enrich Core's capabilities via the GKGS.
- **Skills** are the **specialized motor and processing capabilities**, executing the fine-grained tasks. Their requirements for safe, powerful execution are dependent on Core's instructions and Armory's provisioning of suitable environments, all underpinned by Foundry's sandboxing. Their results are fed back to Core.
- **Guardian** is the **conscience, immune system, and learning center**, ensuring safety, detecting issues, and driving self-improvement. Its requirements for oversight are entirely dependent on the comprehensive observability provided by Foundry and the behavioral data from all other components. Its outputs directly influence Core's learning and enforce constraints on all agents via Foundry's security policies.

The correlations are constant feedback loops: Guardian's evaluation triggers retraining (via Foundry) which improves Skills' models, which in turn leads to better task execution, which Core learns from (via feedback loop with Guardian/GKGS), leading to smarter planning, and so on. Any failure in one dependency cascades, but the redundancy and observability built by Foundry aim to quickly identify and mitigate these.

# Outline

**Categorized Documentation Links:** For each core technology, library, framework, or algorithm, I will provide links to its official documentation, GitHub repositories, and relevant academic or industry standards.

**Conceptual Code Documentation References:** Since I cannot generate an actual, runnable codebase for an AGI, I will detail *what kind* of code documentation (e.g., API definitions, configuration schemas, design patterns, internal module specifications) would be present within a real-world implementation for each component and function. This fulfills the requirement of "no line of code unexamined" from a documentation and architectural design perspective.

**Key Parameters & Configuration Guides:** For major tools, I will explicitly reference sections of their documentation where critical parameters, scaling configurations, and security hardening measures would be found and scrutinized during implementation.

**Best Practices & Design Patterns:** I will point to established best practices and design patterns (e.g., microservices, event-driven architectures, MLOps patterns) that would guide the actual code implementation.

**Technical Stack Overview (Nexus):**

- **Language:** Python (chosen for rapid prototyping and its rich ML ecosystem) and/or Node.js (for high-concurrency I/O in the proxy layer).
- **Frameworks:** FastAPI (Python) for robust API development and Express.js (Node.js) for efficient web applications.
- **Core LLM:** A quantized Mistral-7B/Llama 3-8B model, specifically selected for its balance of performance and efficiency, served by vLLM/Ollama.
- **Database/Cache:** Redis for high-speed caching and session management, and PostgreSQL for long-term, reliable data persistence.
- **Communication:** gRPC for high-performance inter-service communication and Protocol Buffers (Protobuf) for efficient, schema-driven data serialization.
- **Monitoring:** Prometheus for metrics collection, Grafana for visualization and alerting, and OpenTelemetry for distributed tracing.

**Minute Aspects & Technical Dissection:**

**1. User Input & Pre-processing Sub-system**

This sub-system is the absolute first point of contact, meticulously engineered for ingesting raw user input from diverse channels and performing initial, channel-agnostic normalization.

- **1.1. Channel Ingestion Adapters (ChannelAdapter Microservices):**

  - **Purpose:** These adapters fundamentally decouple Nexus's core logic from the intricacies of specific messaging platform APIs. Each dedicated adapter (e.g., SlackAdapter, WebChatAdapter, TelegramAdapter, APIAdapter) is responsible for

translating platform-specific messages into a standardized internal format, ensuring uniformity regardless of the input source.

- ○ **Technical Dissection:**
  - ■ **Standardized Message Schema:** Every incoming message is transformed into a `ScriptKittyMessageSchema` (a Protobuf-defined structure) containing `user_id`, `session_id`, `timestamp`, `channel_type`, `raw_text`, and an `event_type`. This schema ensures strict type-checking and efficient data exchange throughout the system.
  - ■ **Platform-Specific APIs:** Each adapter leverages the official SDKs or REST APIs of its respective platform (e.g., Slack Events API, Telegram Bot API, custom WebSocket connections for web chat).
  - ■ **Event Handling:** Adapters are designed to be event-driven, often utilizing webhooks to receive real-time updates from messaging platforms, minimizing latency.
  - ■ **Input Validation & Sanitization:** Before transformation, inputs are rigorously validated against predefined rules (e.g., maximum message length, allowed character sets) and sanitized to prevent injection attacks or malformed data.
  - ■ **Error Handling & Retries:** Robust mechanisms are in place for handling API rate limits, network failures, and platform-specific errors, often incorporating exponential backoff strategies for retries.
  - ■ **Stateless Design:** Channel adapters are designed to be largely stateless, processing individual messages and forwarding them. Session-specific state is managed by the `SessionManager`.
  - ■ **Deployment:** Each adapter is deployed as a separate, horizontally scalable microservice in Kubernetes, allowing for independent scaling based on the load from specific channels.

- ● **1.2. Input Normalization & Security Filter:**

  - ○ **Purpose:** This component applies a layer of universal pre-processing to the standardized input, enhancing data quality and enforcing initial security checks before NLP processing.
  - ○ **Technical Dissection:**
    - ■ **Unicode Normalization:** All text is converted to a consistent Unicode form (e.g., NFC) to prevent issues with character representation.
    - ■ **Whitespace & Punctuation Cleanup:** Redundant whitespaces are collapsed, and non-standard punctuation is normalized, ensuring cleaner input for downstream NLP.
    - ■ **Basic Spell Correction (Optional/Lightweight):** A lightweight, fast spell-checking library is employed for common typos, improving the robustness of NLU.
    - ■ **PII Redaction (Initial Pass):** Regular expressions or simple pattern matching are used for a preliminary scan and redaction of obvious Personally Identifiable Information (PII) like phone numbers or email

addresses, serving as an initial safety net. More comprehensive PII handling occurs later in the Guardian.

- **Toxicity/Safety Scoring (Heuristic):** A fast, heuristic-based model or rule set provides an initial, low-latency toxicity score, allowing for immediate flagging of highly problematic inputs, preventing them from proceeding deeper into the system. This is not for definitive content moderation but for early detection.
- **Anti-Spam/Flood Control:** Implements mechanisms to detect and mitigate spam or flooding attempts from a single user or IP address, protecting system resources.
- **Dependency:** This filter operates on the `raw_text` field of the `ScriptKittyMessageSchema` received from the Channel Adapters.

**2. Dialogue State & Context Management Sub-system**

This crucial sub-system maintains a persistent, rich conversational context for every user session, enabling coherent, multi-turn interactions.

- **2.1. Session Manager:**

  - **Purpose:** The central orchestrator for session-specific data, including conversational history, identified entities, user preferences, and ongoing task states. It ensures continuity across turns.
  - **Technical Dissection:**
    - **Key-Value Store (Redis):** Utilized for high-speed, low-latency access to active session data. Each `session_id` maps to a JSON blob or a structured object containing `conversation_history` (list of `ScriptKittyMessageSchema` objects), `extracted_entities`, `user_preferences`, `current_task_state`, and `last_activity_timestamp`.
    - **Relational Database (PostgreSQL):** Provides long-term persistence and archival for all sessions, enabling analytical queries and recovery. A robust ORM (e.g., SQLAlchemy) manages the mapping between application objects and database tables.
    - **Data Structure:** The session context is represented as a structured JSON object, allowing for flexible storage of diverse data types and easy extension.
    - **Cache Invalidation & Eviction:** Redis entries have TTL (Time-To-Live) for inactive sessions, and a background process regularly syncs active sessions to PostgreSQL to prevent data loss.
    - **Concurrency Control:** Mechanisms (e.g., optimistic locking) are in place to handle concurrent updates to a session, preventing data corruption.
    - **API Endpoints:** Exposes gRPC endpoints for `get_session_context(session_id)`,

```
update_session_context(session_id, data), and
create_new_session().
```
- **2.2. Contextual Memory & Retrieval:**

  - **Purpose:** Dynamically fetches and updates relevant session data based on the current turn, providing a rich context for the Nexus LLM.
  - **Technical Dissection:**
    - **Conversation History Management:** Manages a sliding window of recent messages within the `conversation_history` array in the session context, ensuring the LLM receives the most relevant past turns without exceeding context limits.
    - **Entity Resolution & Tracking:** Updates extracted entities in the session context, resolving ambiguities and tracking coreferent mentions across turns.
    - **User Preferences & Profiles:** Retrieves and incorporates stored user preferences (e.g., tone, verbosity, domain expertise) into the LLM prompt, personalizing responses.
    - **Active Task State Integration:** If Core initiates a multi-step task, Nexus stores and retrieves the `current_task_state` to guide subsequent user interactions and provide context to Core.
    - **Vector Search (Optional/Future):** For extremely long conversations or external knowledge, a vector database (e.g., Weaviate, Chroma) could store embeddings of past turns or relevant knowledge chunks, allowing semantic retrieval to augment the LLM's context window.

## 3. Core NLU (Natural Language Understanding) & Intent Routing Sub-system

This sophisticated sub-system analyzes user input to identify intent, extract entities, and intelligently route the query to the appropriate internal handler or to Script Kitty.Core.

- **3.1. Intent & Entity Extraction (Nexus LLM):**

  - **Purpose:** The primary NLU engine for general and ambiguous queries, leveraging a fine-tuned LLM to understand nuanced user intent and extract key information.
  - **Technical Dissection:**
    - **LLM Model:** A quantized Mistral-7B/Llama 3-8B model, specifically chosen for its efficient inference and strong NLU capabilities. It's served by `vLLM` or `Ollama` for high-throughput, low-latency inference on GPU/CPU.
    - **System Prompt Engineering:** Carefully crafted system prompts instruct the LLM to act as an NLU engine, specifying the desired JSON output format for intents and entities. This prompt also includes the current `session_context` (from 2.2) to enable contextual understanding.

- - - **Few-Shot Examples:** The prompt may include few-shot examples demonstrating desired intent classification and entity extraction for common scenarios, guiding the LLM's behavior.
      - **Output Schema Validation:** The LLM's JSON output is immediately validated against a `ScriptKittyIntentSchema` (Protobuf) to ensure it conforms to the expected structure before further processing. Retries or fallback mechanisms are in place for invalid outputs.
      - **Confidence Scoring:** The LLM output might include a confidence score for its intent classification, which can be used by the Router for fallback decisions.
      - **Fine-tuning (Continuous):** The Nexus LLM is continuously fine-tuned on anonymized user interactions and human-corrected intent/entity pairs, improving its accuracy and adapting to evolving user language.
- **3.2. Semantic Router & Fallback Mechanism:**

  - **Purpose:** Determines whether Nexus can handle the query directly (e.g., simple chat, small data retrieval) or if it requires escalation to Script Kitty.Core for complex planning and task execution.
  - **Technical Dissection:**
    - **Hybrid Routing Strategy:**
      - **LLM-based Routing (Primary):** The Nexus LLM itself is prompted to classify the intent into broad categories like `GENERAL_CHAT`, `TASK_REQUEST`, `INFORMATION_QUERY`, `FEEDBACK`, etc. It can also suggest which backend agent (Core, Armory, Skills, Guardian) is most relevant if it recognizes specific keywords or patterns.
      - **Traditional Classifier (Fallback/Fast Path):** For high-confidence, well-defined intents (e.g., "hello", "thank you", "what's the weather"), a lightweight, pre-trained transformer-based text classifier (e.g., fine-tuned `bert-base-uncased` from Hugging Face Transformers) provides a faster, more deterministic classification. This acts as a bypass for common queries.
      - **Rule-Based Overrides:** Explicit rules can override LLM or classifier decisions for critical security or routing requirements (e.g., if a "security" keyword is detected, always route to Guardian's policy engine).
    - **Confidence Thresholding:** The router considers confidence scores from both the LLM and the traditional classifier. If confidence is below a certain threshold, the query might be flagged for human review (via Guardian) or routed to Core as a more general request.
    - **Core Escalation:** If the detected intent is `TASK_REQUEST`, `COMPLEX_INFORMATION_QUERY`, or involves capabilities beyond Nexus's direct handling, the standardized

`ScriptKittyMessageSchema` (now enriched with intent and entities) is forwarded to Script Kitty.Core's internal communication hub via gRPC.
- **Direct Nexus Response:** If the intent is `GENERAL_CHAT` or a simple query Nexus can answer directly (e.g., using its internal conversational abilities), it proceeds to the Response Synthesis layer.
- **Orchestration Logic:** Custom Python/Go code orchestrates the decision flow: LLM first, then traditional NLU fallback for known intents, then rule-based overrides, and finally the routing decision.

## 4. Response Synthesis & Persona Layer

This layer transforms structured responses and intermediate updates from Script Kitty.Core and other backend agents into natural, coherent, and persona-consistent chat responses, maintaining the "group chat" illusion.

- **4.1. Response Generation (Nexus LLM):**

  - **Purpose:** Takes the structured output from internal NLU or from Script Kitty.Core's task execution and generates human-like conversational responses.
  - **Technical Dissection:**
    - **LLM Model:** The same Mistral-7B/Llama 3-8B LLM instance used for NLU, but now prompted for text generation.
    - **System Prompt for Persona:** A persistent system prompt ensures the LLM adheres to the "Script Kitty" persona: "empirical, unparalleled inventive, limitlessly creative, unequaled researcher, vast source of information, world-leading expert in Computer Science, AI, Programming, etc." This persona is consistently injected.
    - **Contextual Prompting:** The prompt includes the original user query, the identified intent and entities, the relevant `session_context`, and crucially, the structured `response_data` (e.g., task results, retrieved information) from Core or other agents.
    - **Markdown Formatting:** The LLM is instructed to generate responses using Markdown for readability (e.g., bolding, bullet points, code blocks for technical details).
    - **Guardrails:** Built-in safeguards within the prompting prevent the LLM from generating harmful, off-topic, or non-persona-compliant content.
- **4.2. Group Chat Facilitation & Multi-Agent Attribution:**

  - **Purpose:** Creates the illusion of a "group chat" by attributing responses to specific internal "agents" and managing the flow of information.
  - **Technical Dissection:**
    - **Attribution Prefixes:** When a response originates from a specific backend agent (e.g., Armory, Skills, Guardian), the Nexus prefixes the generated text with clear attributions like "Armory reports:", "Skills has

completed the task:", or "The Nexus has determined:". This is achieved by the Core sending `source_agent` metadata with its responses.

- **Turn Management:** Nexus implicitly manages conversational turns, ensuring responses are coherent and logically sequenced based on the `event_type` and `task_status` signals received from Core.
- **Progress Updates:** For long-running tasks, Nexus receives periodic `PROGRESS_UPDATE` messages from Core and synthesizes them into concise, user-friendly updates (e.g., "Core is currently planning the task...", "Skills is executing the code...").
- **Templating Engine:** For highly structured or repetitive responses (e.g., error messages, status reports), a templating engine (like Jinja2 in Python) can be used to insert dynamic data into predefined conversational templates, ensuring consistency.

**5. Intermediary & Internal Communication Proxy**

This sub-system serves as the gatekeeper for all external user-facing communication and the forwarding hub for structured requests to Script Kitty.Core.

- **5.1. Core Communication Gateway (gRPC Client):**

  - **Purpose:** The dedicated conduit for all structured communication between Nexus and Script Kitty.Core.
  - **Technical Dissection:**
    - **gRPC Client Implementation:** Nexus runs a gRPC client that connects to Script Kitty.Core's gRPC server. This provides high-performance, bidirectional streaming capabilities for complex interactions.
    - **Protobuf Message Exchange:** All messages (`ScriptKittyMessageSchema`, `TaskRequestSchema`, `TaskUpdateSchema`, `TaskResultSchema`) are defined using Protocol Buffers, ensuring type safety, backward/forward compatibility, and efficient serialization/deserialization across different services and languages.
    - **Asynchronous Communication:** While gRPC can be synchronous, the client is designed to handle asynchronous responses from Core, particularly for long-running tasks where Core might send multiple `TaskUpdate` messages before a final `TaskResult`.
    - **Connection Management:** Robust connection pooling and retry logic for maintaining a stable connection to Core.
- **5.2. User Output Renderer & Channel Dispatcher:**

  - **Purpose:** Takes the final, persona-infused conversational response and dispatches it back to the user via the appropriate channel.
  - **Technical Dissection:**

- **Standardized Output Schema:** Responses are encapsulated in a `ScriptKittyOutputSchema` (Protobuf), containing the `session_id`, `channel_type`, `response_text` (Markdown), and `suggested_actions` (e.g., buttons, quick replies).
- **Channel Adapter Integration:** Reuses the same `ChannelAdapter` microservices from 1.1 (e.g., SlackAdapter, WebChatAdapter) but in reverse. The adapter translates the `ScriptKittyOutputSchema` into the platform-specific message format (e.g., Slack blocks, Telegram keyboard).
- **Error Handling:** Catches and logs errors during message dispatch (e.g., network issues, platform API errors), potentially triggering retries or fallback notifications.
- **Response Timing & Throttling:** Ensures responses are sent at an appropriate pace, preventing overwhelming the user or hitting channel-specific rate limits.
- **Interactive Components:** If `suggested_actions` are included, the dispatcher handles their rendering as interactive elements (e.g., buttons in Slack, web chat).

**Operational Excellence & Scalability (Nexus-Specific):**

- **Observability (Integrated with Foundry):**

  - **Metrics:** Prometheus exporters collect crucial metrics like `requests_per_second`, `latency_per_endpoint`, `llm_inference_time`, `active_sessions_count`, and `error_rates`.
  - **Logging:** Centralized logging to Loki/Elasticsearch via Fluent Bit for all Nexus microservices, categorized by service, session ID, and log level, facilitating rapid debugging.
  - **Tracing:** OpenTelemetry SDKs embedded in all Nexus services provide distributed traces, allowing end-to-end visibility of user requests across multiple internal components.
  - **Monitoring & Alerting:** Grafana dashboards visualize these metrics in real-time. Prometheus/Grafana alert managers trigger critical alerts (e.g., `Nexus_LLM_Latency_High`, `Nexus_Core_Connection_Errors`) to on-call teams via PagerDuty/Slack.
- **Scalability:**

  - **Stateless Channel Adapters:** Designed to be highly scalable via Kubernetes Horizontal Pod Autoscaling (HPA), reacting to increased message queue depth or active connections.
  - **SessionManager:** Utilizes a Redis cluster for horizontal scaling of caching and PostgreSQL replication/sharding for robust long-term persistence.

- ○ **LLM Inference:** The dedicated vLLM/Ollama service supports horizontal scaling of GPU/CPU instances, efficiently distributing inference load.
  - ○ **Kubernetes HPA:** Automatically scales Nexus pods based on CPU/memory utilization or custom metrics (e.g., `messages_in_queue_depth` or `active_websocket_connections`).
  - ○ **Service Mesh Load Balancing:** Istio or Linkerd automatically load balances requests across healthy pods within each Nexus service, ensuring even distribution of traffic and fault tolerance.
- **Continuous Evolution (Via Script Kitty.Foundry):**

  - ○ **Containerization:** Every Nexus microservice is rigorously packaged into immutable Docker containers, ensuring consistency and portability across environments.
  - ○ **CI/CD Pipelines:** Automated build, test, and deployment pipelines (e.g., using GitHub Actions or GitLab CI/CD) are triggered by Git commits to Nexus's code repositories, ensuring rapid and reliable delivery.
  - ○ **Seamless Deployment:** New versions of Nexus services or updates to the Nexus LLM model are deployed using Blue/Green or Canary deployment strategies, orchestrated by Kubernetes and Argo CD/Flux CD. This ensures zero downtime during updates and allows for rapid rollbacks if issues are detected post-deployment.
  - ○ **Model Updates:** The Nexus LLM itself undergoes continuous fine-tuning based on anonymized user interactions and human feedback. New model artifacts are packaged into fresh container images for the vLLM service, which are then seamlessly swapped out by Kubernetes, giving the impression of a continuously improving conversational AI.

---

## II. Script Kitty.Core - The Central Governing AI / High-Level Orchestrator

**Core Purpose:** Script Kitty.Core is the strategic brain of the entire Script Kitty system. It receives high-level goals from Script Kitty.Nexus, breaks them down into executable plans, orchestrates the various specialized backend AIs (Armory, Skills, Guardian), manages global system state, and integrates feedback to drive overall system intelligence. It embodies the "overarching capacity, functionality, and knowledge reinforcement" of Script Kitty.

**Code Summary with Title:**

- **`CoreService/main.py`**: The entry point for the Core microservice, setting up gRPC server and initializing planning and orchestration engines.
- **`CoreService/planning_engine.py`**: Contains the HTN (Hierarchical Task Network) planning logic, LLM integration for high-level planning, and symbolic planning fallback mechanisms.

- **`CoreService/orchestrator.py`**: Manages the lifecycle of specialized agents (Armory, Skills, Guardian) via Kubernetes API, assigns tasks, and monitors progress.
- **`CoreService/knowledge_graph.py`**: Implements the Global Knowledge Graph (GKGS) interface for storing and retrieving facts, rules, and task contexts, potentially integrating with Neo4j/ArangoDB clients and Weaviate/Chroma for vector search.
- **`CoreService/message_broker.py`**: Manages internal communication using gRPC for direct calls and Kafka/Pulsar for asynchronous event streaming between Core and other agents.
- **`CoreService/feedback_loop.py`**: Ingests metrics, success/failure signals, and human corrections, feeding them into the GKGS and triggering learning mechanisms.
- **`CoreService/schemas/core_messages.proto`**: Protobuf definitions for `PlanSchema`, `TaskAssignment`, `AgentStatus`, and other messages exchanged between Core and other agents.
- **`CoreService/data/plan_templates.json`**: A repository of pre-defined plan templates used by the planning engine for common tasks.
- **`CoreService/config.yaml`**: Configuration file for Core's operational parameters, including LLM endpoints, database connection strings, and agent registry mappings.

**Minute Aspects & Technical Dissection: [Placeholder for detailed analysis of Script Kitty.Core]**

---

## III. Script Kitty.Armory - Resource Acquisition & Tooling AI

**Core Purpose:** Script Kitty.Armory is the system's external sensory and manipulation apparatus, specializing in acquiring, processing, and digesting external information (web scraping, research), managing a dynamic tool registry, and providing secure and controlled access to external resources (e.g., cloud services, specialized APIs, quantum computing interfaces). It is the source of Script Kitty's vast and ever-expanding actionable intelligence.

**Code Summary with Title:**

- **`ArmoryService/main.py`**: The main entry point for the Armory microservice, configuring API endpoints and integrating with sub-components.
- **`ArmoryService/web_scraper.py`**: Implements intelligent web navigation, data extraction (using Playwright/Selenium, BeautifulSoup4), and content digestion (Trafilatura, LLM summarization).
- **`ArmoryService/nlu_extractor.py`**: Contains NLP models (spaCy, Hugging Face Transformers) for named entity recognition (NER) and relation extraction (RE) from scraped text.
- **`ArmoryService/tool_registry_api.py`**: Exposes API endpoints for Core to query and manage the dynamic tool registry, backed by a PostgreSQL database.

- **`ArmoryService/tool_manifest_parser.py`**: Logic for parsing tool documentation (e.g., OpenAPI specs, CLI man pages) to generate structured `ToolManifestSchema` entries.
- **`ArmoryService/wrapper_generator.py`**: Uses a code-focused LLM (CodeLlama/StarCoder2) and templating (Jinja2) to generate secure API/CLI wrappers, with integration points for human review.
- **`ArmoryService/iac_provisioner.py`**: Interfaces with Terraform, Ansible, and cloud provider SDKs to provision computational resources based on Core's requests. Includes quantum computing SDK integrations (Qiskit, Cirq).
- **`ArmoryService/schemas/armory_messages.proto`**: Protobuf definitions for `ResearchQuery`, `ToolManifest`, `ResourceRequest`, `GeneratedWrapperCode`.
- **`ArmoryService/templates/wrapper_templates/`**: Directory containing secure code templates for various programming languages to be filled by the wrapper generator.
- **`ArmoryService/policies/resource_provisioning.rego`**: OPA policies for governing resource provisioning requests (integrated with Guardian).

**Minute Aspects & Technical Dissection: [Placeholder for detailed analysis of Script Kitty.Armory]**

---

## IV. Script Kitty.Skills - Task Execution AI

**Core Purpose:** Script Kitty.Skills houses the specialized AI models and agents that execute concrete computational tasks identified by Script Kitty.Core. These are the "doers" of the system, performing low-level data analysis, executing generated code, processing images, generating creative content, or solving specific algorithmic problems within highly secure and isolated environments.

**Code Summary with Title:**

- **`SkillsService/main.py`**: The entry point for the Skills microservice, managing task queues and dispatching to specialized executors.
- **`SkillsService/model_inference_api.py`**: Exposes gRPC endpoints for various specialized ML models (e.g., computer vision, time-series, code generation) served by TorchServe/TensorFlow Serving/KServe.
- **`SkillsService/sandbox_executor.py`**: Manages the secure, containerized execution environment (Docker, Firejail, seccomp-bpf) for arbitrary code (Python, R) and CLI commands.
- **`SkillsService/data_processor.py`**: Contains modules for data manipulation, cleaning, and statistical analysis using libraries like Pandas, NumPy, Dask, Polars.
- **`SkillsService/generative_executor.py`**: Orchestrates code-focused LLMs (CodeLlama, DeepSeek Coder) for generating code snippets, unit tests, or structured reports based on task instructions. Includes AST parsers for code validation.

- **`SkillsService/problem_solver_modules/`**: Directory containing highly specialized problem-solving algorithms (e.g., optimization algorithms, symbolic solvers, domain-specific rule engines).
- **`SkillsService/schemas/skills_messages.proto`**: Protobuf definitions for `TaskExecutionRequest`, `CodeExecutionResult`, `ModelInferenceResult`, `TaskProgressUpdate`.
- **`SkillsService/sandbox_configs/`**: YAML files defining Docker container parameters, resource limits, and security profiles (seccomp, AppArmor) for sandboxed environments.
- **`SkillsService/tests/generated_code_tests.py`**: Unit tests for validating the functionality of code generated by the `generative_executor`.

**Minute Aspects & Technical Dissection: [Placeholder for detailed analysis of Script Kitty.Skills]**

---

# V. Script Kitty.Guardian - Safety & Alignment Proxy / Evaluation & Training AI

**Core Purpose:** Script Kitty.Guardian is the system's conscience, its ethical compass, and its continuous improvement engine. It vigilantly enforces ethical, legal, and safety controls, monitors system performance, identifies learning opportunities, and triggers training/retraining cycles for all Script Kitty models, ensuring alignment with its core principles and human values.

**Code Summary with Title:**

- **`GuardianService/main.py`**: The central Guardian microservice, orchestrating policy enforcement, monitoring, and learning loops.
- **`GuardianService/policy_engine.py`**: Integrates Open Policy Agent (OPA) for declarative policy enforcement (Rego rules) and potentially a custom rule engine for complex ethical dilemmas.
- **`GuardianService/llm_safety_evaluator.py`**: Utilizes a fine-tuned Ethical LLM to identify potential ethical conflicts within proposed actions, integrated with XAI tools (LIME, SHAP).
- **`GuardianService/monitoring_agent.py`**: Subscribes to metrics (Prometheus) and logs (Loki/Elasticsearch) streams, performing real-time anomaly detection (PyOD) on system performance and behavior.
- **`GuardianService/evaluation_trigger.py`**: Based on monitoring data, human feedback, or schedule, triggers automated evaluation and retraining pipelines via Kubeflow Pipelines/Argo Workflows.
- **`GuardianService/bias_detector.py`**: Analyzes training data and model outputs for biases (using Aequitas, Fairlearn, AI Fairness 360) and provides mitigation strategies.

- **GuardianService/feedback_api.py**: Provides API endpoints for human operators to provide explicit approval, override, or corrective feedback via a dedicated UI.
- **GuardianService/schemas/guardian_messages.proto**: Protobuf definitions for `ActionProposal`, `PolicyViolation`, `SystemMetric`, `FeedbackEntry`.
- **GuardianService/policies/security_policies.rego**: OPA policies for system-wide security, access control, and data handling.
- **GuardianService/policies/ethical_guidelines.txt**: Text-based guidelines for the Ethical LLM and human review.
- **GuardianService/ml_pipeline_definitions/**: YAML/Python definitions for Kubeflow/Argo Workflows that manage model training and evaluation.

**Minute Aspects & Technical Dissection: [Placeholder for detailed analysis of Script Kitty.Guardian]**

---

# VI. Script Kitty.Foundry - The MLOps & Infrastructure Fabric

**Core Purpose:** Script Kitty.Foundry is the bedrock, the unseen yet indispensable engine that propels the entire Script Kitty ecosystem. It is the sophisticated MLOps (Machine Learning Operations) and infrastructure orchestration layer, providing universal services for compute, data, security, deployment, and operational intelligence. Foundry ensures seamless transitions, continuous evolution, universal agent features, and a platform that learns and adapts with the collective knowledge of all Script Kitty agents, making the abstract notion of "AGI" a robust and practical reality. It underpins every single function, component, and interaction within Script Kitty, guaranteeing its unparalleled functionality and relentless improvement.

**Technical Stack Overview (Foundry):**

- **Orchestration:** Kubernetes (the foundational distributed operating system for Script Kitty).
- **Containerization:** Docker.
- **CI/CD & GitOps:** GitLab CI/CD/GitHub Actions (CI/CD Pipelines), Argo CD/Flux CD (GitOps).
- **Observability:** Prometheus (metrics), Grafana (dashboards, alerts), Loki/Elasticsearch (logging), Fluent Bit/Fluentd (log collection), OpenTelemetry (tracing), Jaeger/Zipkin (trace visualization).
- **Data Management:** MinIO/Ceph/AWS S3/GCP GCS (Object Storage), DVC (Data Version Control), Feast (Feature Store), Apache Airflow/Prefect/Dagster (Data Orchestration).
- **Security:** HashiCorp Vault (secrets), Calico/Cilium (network policy), Istio/Linkerd (service mesh), Falco (runtime security), Trivy/Clair (container scanning), Keycloak (identity management).
- **ML Pipelines:** Kubeflow Pipelines/Argo Workflows (workflow orchestration), MLflow (experiment tracking), Ray Tune/Optuna (hyperparameter optimization).

- **Virtualization (Hybrid):** KubeVirt (for specialized legacy workloads if needed).

**Minute Aspects & Technical Dissection:**

**1. Kubernetes-Native Compute & Orchestration Sub-system**

This is the very operating system of Script Kitty's distributed brain. It is responsible for the deployment, scaling, healing, and resource management of every microservice and specialized agent.

- **1.1. Kubernetes Cluster Management (Control Plane & Worker Nodes):**

  - **Purpose:** Provides the declarative, self-healing, and scalable foundation for all Script Kitty operations.
  - **Technical Dissection:**
    - **API Server:** The central control plane component, exposing the Kubernetes API (RESTful interface) for all interactions. Every deployment, scaling event, or configuration change for Script Kitty's agents goes through this API. Authentication and authorization (via RBAC - Role-Based Access Control) are meticulously enforced for all API calls.
    - **etcd:** The highly available key-value store, serving as Kubernetes's consistent and persistent backbone for cluster state and configuration. All desired states for Script Kitty's services are stored here.
    - **Scheduler:** Monitors newly created pods and assigns them to worker nodes based on resource requirements (CPU, memory, GPU), node constraints (taints and tolerations), and affinity/anti-affinity rules. Foundry's scheduler ensures optimal placement of Script Kitty's diverse agents (e.g., placing GPU-intensive LLMs on GPU nodes).
    - **Controller Manager:** Runs various controllers (e.g., Deployment Controller, ReplicaSet Controller) that continuously watch the state of the cluster and make changes to move the current state towards the desired state. For Script Kitty, this means ensuring the correct number of Nexus, Core, Armory, Skills, and Guardian pods are always running.
    - **Kubelet (on Worker Nodes):** The agent that runs on each worker node, ensuring containers are running in a Pod. It communicates with the API server, manages Pods, mounts volumes, and reports node status. Kubelet is critical for executing Script Kitty's agent workloads.
    - **Kube-proxy (on Worker Nodes):** Maintains network rules on nodes, enabling network communication to your Pods from inside or outside the cluster. It ensures that Script Kitty's internal microservices can communicate seamlessly and reliably.
    - **Horizontal Pod Autoscaler (HPA):** Dynamically scales the number of running pods for stateless services (e.g., Nexus Channel Adapters, LLM inference services, Core's message processors) based on observed CPU utilization, memory consumption, or custom metrics (e.g., inbound

message queue depth, active user sessions). This ensures Script Kitty can handle fluctuating loads without manual intervention.

- ■ **Vertical Pod Autoscaler (VPA) (Experimental/Advisory):** While primarily used for recommendations, VPA can adjust resource requests/limits for pods, optimizing resource allocation based on historical usage. This helps right-size Script Kitty's specialized agents over time.
- ■ **Cluster Autoscaler:** Automatically adjusts the number of worker nodes in the Kubernetes cluster based on pending pods and resource utilization. This is critical for scaling Script Kitty's entire compute fabric up and down dynamically with demand, optimizing cloud costs.
- ■ **Node Affinity/Anti-Affinity & Taints/Tolerations:** Used extensively to ensure specific agent types (e.g., GPU-intensive LLMs for Core/Skills, high-I/O data processing for Armory) are scheduled on appropriate nodes, or to prevent co-location of sensitive services on the same node for security/performance reasons.

- ● **1.2. Hybrid Compute Capabilities (KubeVirt Integration - Conditional):**

  - ○ **Purpose:** While primarily container-native, Foundry can optionally integrate KubeVirt to run virtual machines alongside containers within the same Kubernetes cluster. This provides flexibility for incorporating legacy tools or specific software that absolutely requires a VM environment (e.g., highly specialized proprietary quantum computing simulators, older enterprise software).
  - ○ **Technical Dissection:**
    - ■ **KubeVirt CRDs:** Custom Resource Definitions (CRDs) like `VirtualMachine` are installed, allowing VMs to be managed like any other Kubernetes resource.
    - ■ **VM Lifecycle Management:** KubeVirt controllers manage the lifecycle of VMs (start, stop, pause, migrate) using standard Kubernetes commands.
    - ■ **Network Integration:** VMs share the same Kubernetes network fabric, allowing seamless communication with containerized services.
    - ■ **Persistent Volumes:** VMs can utilize Kubernetes Persistent Volumes, integrating with Foundry's storage management.
    - ■ **Use Cases for Script Kitty:** Potentially for specialized "Skills" agents that rely on non-containerizable software, or for a highly isolated execution environment for extremely sensitive "Armory" tool wrappers.

- ● **1.3. Custom Resource Definitions (CRDs) for Script Kitty Agents:**

  - ○ **Purpose:** Foundry defines custom Kubernetes resources that represent Script Kitty's specialized agents (e.g., `ScriptKittyAgent`, `SkillModule`, `ToolManifest`). This allows Kubernetes to understand and manage these high-level logical constructs directly.
  - ○ **Technical Dissection:**

- **YAML Definitions:** CRDs are defined in YAML files, specifying their schema (`spec` and `status`) and versioning.
- **Custom Controllers:** Foundry implements custom Kubernetes controllers (often using operator frameworks like Kubebuilder or Operator SDK) that watch for changes to these CRDs. When a `ScriptKittyAgent` CR is created or updated, the custom controller translates this into standard Kubernetes Deployments, Services, and other resources.
- **Declarative Agent Management:** This enables Core or even human operators to declaratively define desired agent configurations (e.g., "deploy an instance of CodeLlama-7B skill module with 2 GPUs") via CRs, and Kubernetes handles the underlying orchestration.
- **Status Reporting:** The custom controller updates the `status` field of the CRD, reflecting the real-time operational status of the agent (e.g., `Running`, `Degraded`, `Updating`).

## 2. Container Registry & Image Management Sub-system

The secure repository for all of Script Kitty's executable components, ensuring versioning, integrity, and controlled distribution.

- **2.1. Secure Container Registry:**
  - **Purpose:** Stores all Docker images for Script Kitty components (Nexus, Core, Guardian, Foundry services) and agent models (Armory's tool wrappers, Skills' specialized ML models).
  - **Technical Dissection:**
    - **Choice of Registry:**
      - **Self-hosted:** Harbor or Quay.io for complete control, enhanced security features (Vulnerability Scanning, Content Trust, Replication), and integration with internal networks.
      - **Cloud-Managed:** AWS ECR, GCP GCR, Azure Container Registry for ease of management, high availability, and integration with cloud-native security services.
    - **Image Tagging & Versioning:** Strict adherence to semantic versioning (e.g., `v1.2.3`, `v1.2.3-commitsha`, `latest-stable`) for all images, facilitating rollbacks and ensuring traceability.
    - **Access Control:** Robust RBAC (Role-Based Access Control) integrated with Script Kitty's central identity management (Keycloak) to strictly control who can push, pull, or delete images.
    - **Image Signing (Content Trust):** Utilizes Notary (or equivalent) for cryptographically signing container images, verifying their authenticity and integrity before deployment to prevent tampering or supply chain attacks.
    - **Vulnerability Scanning:** Integrated scanners (e.g., Trivy, Clair, Harbor's built-in scanner) automatically scan all pushed images for known CVEs

(Common Vulnerabilities and Exposures), flagging insecure dependencies and preventing deployment of vulnerable components.
- **Image Immutability:** Once an image is tagged and pushed, it is immutable, ensuring that the deployed artifact is precisely what was built and tested.

**3. CI/CD & GitOps Automation Sub-system**

The core of Script Kitty's continuous evolution, automating the entire software delivery lifecycle with a strong emphasis on declarative configuration and version control.

- **3.1. Continuous Integration (CI) Pipelines:**

  - **Purpose:** Automates the build and test phases for every code change made across Script Kitty's repositories.
  - **Technical Dissection:**
    - **Triggering:** Pipelines are automatically triggered by Git events (e.g., `push` to `main` branch, `pull_request` creation).
    - **Build Jobs:** Compiles source code (e.g., Python, Go), resolves dependencies, and builds Docker images for microservices and agent models.
    - **Unit Tests:** Executes comprehensive unit tests for individual functions and classes.
    - **Integration Tests:** Validates interactions between closely related components (e.g., Nexus NLU with Session Manager).
    - **Linting & Static Analysis:** Enforces coding standards and detects potential bugs or security vulnerabilities (e.g., `Pylint`, `ESLint`, `Bandit`, `SonarQube` community edition).
    - **Security Scans (Pre-Push):** Runs basic vulnerability scans on dependencies and code before image creation.
    - **Artifact Generation:** If all tests pass, it generates deployable artifacts (e.g., Docker images, Helm charts) and pushes them to the secure container registry with appropriate version tags.
    - **Tools:** GitLab CI/CD, GitHub Actions, Jenkins, Tekton Pipelines are all viable orchestrators for these complex, multi-stage pipelines.
- **3.2. Continuous Delivery/Deployment (CD) & GitOps:**

  - **Purpose:** Automates the deployment and update processes for all Script Kitty services and agents, ensuring that the live system always reflects the desired state defined in Git.
  - **Technical Dissection:**
    - **Declarative Infrastructure as Code (IaC):** All Kubernetes manifests (Deployments, Services, ConfigMaps, CRDs, Helm Charts) are stored in

a Git repository, serving as the single source of truth for Script Kitty's infrastructure and application configurations.

- **GitOps Agents:** Argo CD or Flux CD are deployed within the Kubernetes cluster. These agents continuously monitor the Git repositories for changes.
- **Automated Synchronization:** When a change is detected in Git (e.g., a new Docker image tag for Nexus), the GitOps agent automatically pulls the new configuration and applies it to the Kubernetes cluster, reconciling the actual state with the desired state.
- **Deployment Strategies:**
  - **Rolling Updates (Default):** Kubernetes natively supports rolling updates, gradually replacing old pods with new ones.
  - **Blue/Green Deployments:** For critical services, Foundry orchestrates Blue/Green deployments (via GitOps or custom controllers). A new "Green" environment is deployed alongside the existing "Blue" one. Once verified, traffic is instantly shifted to Green, minimizing downtime. Rapid rollback to Blue is possible.
  - **Canary Deployments:** A small percentage of user traffic is routed to new "Canary" version of a service. Foundry's monitoring (via Guardian's integration) meticulously tracks metrics (latency, error rates, user feedback) for the Canary. If performance is stable, traffic is gradually shifted; otherwise, the Canary is rolled back. This is crucial for safely deploying new LLM models or complex agent logic.
- **Automated Rollbacks:** If a deployment fails (e.g., pods crash, health checks fail, or Guardian detects anomalies post-deployment), the GitOps agent automatically rolls back to the last known good configuration in Git.
- **Immutable Deployments:** Once a container image is built and tested, it's deployed as is. Configuration changes are applied via ConfigMaps or Secrets, which trigger rolling updates.

## 4. Observability Stack (Centralized Logging, Metrics, Tracing) Sub-system

The eyes and ears of Script Kitty, providing real-time insights into its operational health, performance, and behavior. Essential for debugging, monitoring, and enabling continuous improvement.

- **4.1. Metrics Collection & Monitoring (Prometheus & Grafana):**

  - **Purpose:** Gathers high-frequency numerical data about the system's state and performance.
  - **Technical Dissection:**
    - **Prometheus:** A powerful time-series database and monitoring system. Each Script Kitty microservice (Nexus, Core, Armory, Skills, Guardian, Foundry components) exposes a `/metrics` endpoint in the Prometheus

text format, containing granular metrics (e.g., `http_requests_total`, `api_call_latency_seconds`, `llm_inference_duration_milliseconds`, `gpu_utilization_percent`, `agent_task_success_rate`).

- **Prometheus Operator:** Deployed in Kubernetes, it automatically discovers and scrapes metrics from all Script Kitty services based on Kubernetes labels/annotations.
- **Grafana:** A visualization and analytics platform. Pre-built and custom dashboards provide real-time views of Script Kitty's health (e.g., overall system latency, individual agent error rates, resource saturation, active user sessions, task completion rates).
- **PromQL:** Complex queries are written using PromQL (Prometheus Query Language) to analyze trends, aggregate data, and identify specific performance bottlenecks within Script Kitty.
- **Alerting:** Prometheus Alertmanager integrates with Grafana to trigger critical alerts (e.g., "Core planning failures exceeding threshold," "Skills sandboxed execution timeout," "Guardian policy violations detected") via Slack, PagerDuty, or email.

- **4.2. Centralized Logging (Fluent Bit/Fluentd & Loki/Elasticsearch):**

  - **Purpose:** Aggregates and stores all log messages generated by Script Kitty's components, enabling debugging, auditing, and post-mortem analysis.
  - **Technical Dissection:**
    - **Log Collection Agents:** Fluent Bit (lightweight, high-performance) or Fluentd are deployed as DaemonSets on each Kubernetes node. They collect container logs (from stdout/stderr) and node-level logs.
    - **Log Aggregation:**
      - **Loki:** Preferred for its Prometheus-inspired query language (LogQL) and efficient storage of semi-structured logs. Ideal for troubleshooting, enabling fast searches for specific `session_id`, `agent_name`, or `error_message`.
      - **Elasticsearch/OpenSearch:** For more advanced full-text search, complex aggregations, and long-term archival of highly structured logs.
    - **Log Enrichment:** Logs are automatically enriched with Kubernetes metadata (pod name, namespace, container ID) and Script Kitty-specific identifiers (session ID, task ID) for easier filtering and correlation.
    - **Log Retention Policies:** Configurable policies ensure efficient storage management, balancing compliance requirements with cost (e.g., hot storage for recent logs, cold storage for older logs).
- **4.3. Distributed Tracing (OpenTelemetry & Jaeger/Zipkin):**

- **Purpose:** Provides end-to-end visibility into the flow of a single request or task execution across multiple microservices within Script Kitty. Crucial for understanding latency, identifying bottlenecks, and debugging complex distributed interactions.
- **Technical Dissection:**
  - **OpenTelemetry SDKs:** Integrated into the code of every Script Kitty microservice (Nexus, Core, Armory, Skills, Guardian, Foundry's own services). These SDKs automatically instrument incoming/outgoing requests, database calls, and function executions.
  - **Span Generation:** Each operation within a request (e.g., "Nexus receives query," "Core plans task," "Armory fetches tool," "Skills executes code") generates a "span." Spans contain metadata (start/end time, duration, attributes like `http.method`, `db.statement`, `agent.name`, `task.id`).
  - **Trace Context Propagation:** Unique `trace_id` and `span_id` are propagated across service boundaries (e.g., HTTP headers, gRPC metadata), linking all related spans into a single "trace."
  - **Trace Exporters:** Spans are exported to a backend collector (e.g., OpenTelemetry Collector).
  - **Jaeger/Zipkin:** Backend systems for storing, indexing, and visualizing traces. Developers and SREs can use Jaeger UI to visualize a flame graph of a single user request, showing which services were involved, how long each step took, and where errors occurred.
  - **Correlation with Logs/Metrics:** Trace IDs are often included in log messages, allowing direct navigation from a trace view to relevant log entries.

**5. Data Management (Data Lake, Data Versioning, Feature Store) Sub-system**

The robust infrastructure for all data assets within Script Kitty, from raw input to trained models, ensuring integrity, accessibility, and reproducibility.

- **5.1. Scalable Object Storage (Data Lake Foundation):**

  - **Purpose:** Provides durable, cost-effective, and massively scalable storage for all unstructured and semi-structured data generated or consumed by Script Kitty.
  - **Technical Dissection:**
    - **Choice of Storage:**
      - **Self-hosted:** MinIO (S3-compatible object storage, can run on Kubernetes) or Ceph (distributed object, block, and file storage for large-scale deployments).
      - **Cloud-Managed:** AWS S3, GCP GCS, Azure Blob Storage (for unparalleled scalability, durability, and integration with other cloud services).

- **Data Types:** Stores raw user inputs, web scraped data, extracted research documents, model training datasets, pre-processed features, model checkpoints, inference logs, audit logs, and Guardian's feedback data.
- **Bucket/Prefix Organization:** Data is logically organized into buckets (or prefixes within a single bucket) based on data source, type, and stage of processing (e.g., `raw-user-inputs`, `armory-research-docs`, `skills-model-checkpoints`, `guardian-feedback-loop`).
- **Lifecycle Policies:** Configured for automatic tiering (e.g., moving older data from hot to cold storage) and deletion, optimizing storage costs.
- **Access Control:** Granular IAM policies control read/write access to specific buckets or prefixes for different Script Kitty agents.

- **5.2. Data Version Control (DVC - Data Version Control):**

  - **Purpose:** Versions datasets and ML models like source code, ensuring reproducibility of experiments and deployments.
  - **Technical Dissection:**
    - **Git Integration:** DVC works alongside Git. Git tracks small `.dvc` files that point to the actual data files, which are stored in object storage.
    - **Data Checksums:** DVC generates checksums for data files, allowing it to detect changes.
    - **Pipeline Definition:** DVC can define data processing and ML training pipelines (`dvc.yaml`), linking code, data, and models. This is critical for Script Kitty's continuous learning. For example, Guardian triggering a retraining of a Skills agent model will use DVC to fetch the exact version of the training data and code.
    - **Reproducibility:** Any team member can `dvc checkout` a specific version of the data and code from Git, ensuring they can reproduce exact training runs or model evaluations.

- **5.3. Feature Store (Feast):**

  - **Purpose:** Provides a centralized, standardized, and discoverable repository for machine learning features, ensuring consistency between training and inference environments.
  - **Technical Dissection:**
    - **Offline Store:** Typically uses a data lake (S3/GCS) or a data warehouse (Snowflake/BigQuery) for historical feature data used in model training by Guardian.
    - **Online Store:** Uses a low-latency key-value store (e.g., Redis, Cassandra) for real-time feature serving during model inference by Skills agents (e.g., fetching a user's recent activity score for a personalized response, or a tool's success rate for Armory's decision making).

- **Feature Definition:** Features are defined declaratively in `feature_store.yaml` files, including data types, sources, and transformation logic.
- **Point-in-Time Correctness:** Feast ensures that historical features used for training accurately reflect the state of the world at the time the events occurred, preventing data leakage.
- **Use Cases for Script Kitty:**
    - **Skills:** Provides pre-computed features for specialized AI models (e.g., historical user interaction patterns for a recommendation model).
    - **Guardian:** Stores performance metrics and feedback signals as features for models that predict optimal retraining schedules or detect anomalies.
    - **Core:** Might use features related to past task success rates for its planning engine.
- **5.4. Data Orchestration (Apache Airflow/Prefect/Dagster):**

    - **Purpose:** Manages and schedules complex data pipelines for ingestion, transformation, and movement of data within Script Kitty's data lake.
    - **Technical Dissection:**
        - **DAGs (Directed Acyclic Graphs):** Data pipelines are defined as DAGs in Python code, representing sequences of tasks with dependencies.
        - **Scheduling:** Tasks can be scheduled at fixed intervals (e.g., daily ingestion of new web data for Armory) or triggered by events (e.g., new model artifact pushed).
        - **Error Handling & Retries:** Robust mechanisms for handling task failures, with configurable retries, alerting, and manual intervention points.
        - **Monitoring & UI:** Provides a web UI for monitoring pipeline runs, viewing logs, and manually triggering tasks.
        - **Use Cases for Script Kitty:**
            - **Armory:** Orchestrates web scraping, data cleaning, and data loading into the data lake.
            - **Guardian:** Manages pipelines for collecting feedback data, generating evaluation metrics, and preparing datasets for model retraining.
            - **Foundry's Internal Data:** Processes logs and metrics for long-term analytics.

## 6. Security & Identity Management Sub-system

The paramount layer ensuring the integrity, confidentiality, and availability of Script Kitty, protecting it from external threats and internal misuse.

- **6.1. Secrets Management (HashiCorp Vault):**

- ○ **Purpose:** Securely stores, accesses, and manages sensitive credentials (API keys, database passwords, cloud tokens) required by Script Kitty's services and agents.
- ○ **Technical Dissection:**
  - ■ **Centralized Storage:** All secrets are stored in an encrypted, audited, and versioned central repository.
  - ■ **Dynamic Secrets:** Vault can generate on-demand, short-lived credentials for databases, cloud providers, and other services, minimizing the risk of long-lived, static secrets. For example, a Skills agent needing to access a database will request a temporary credential from Vault.
  - ■ **Leasing & Revocation:** Secrets have associated leases and can be automatically revoked after use or expiry.
  - ■ **Audit Logging:** All access to secrets is meticulously logged for auditing and compliance.
  - ■ **Authentication & Authorization:** Integrates with Kubernetes Service Accounts and other identity providers (e.g., Keycloak) to authenticate clients requesting secrets and authorize access based on policies.
- ● **6.2. Network Policy (Calico/Cilium):**

  - ○ **Purpose:** Enforces network segmentation and traffic filtering within the Kubernetes cluster, creating a "zero-trust" environment where agents can only communicate with explicitly allowed services.
  - ○ **Technical Dissection:**
    - ■ **Declarative Rules:** Network policies are defined as Kubernetes YAML manifests, specifying allowed ingress (inbound) and egress (outbound) traffic based on labels, namespaces, IP blocks, and ports.
    - ■ **Micro-segmentation:** Enforces isolation between Script Kitty's services (e.g., Nexus can talk to Core, but Skills cannot directly talk to Core's sensitive internal endpoints unless explicitly allowed).
    - ■ **Prevention of Lateral Movement:** Restricts an attacker from moving laterally across the cluster if one service is compromised. For example, a compromised Skills agent cannot initiate connections to the Guardian's policy enforcement engine or Foundry's secrets management.
    - ■ **DNS Policies:** Can restrict DNS lookups to authorized internal services only.
- ● **6.3. Service Mesh (Istio/Linkerd - Security Features):**

  - ○ **Purpose:** Adds a layer of sophisticated traffic management, observability, and *enhanced security* capabilities to Script Kitty's inter-service communication without modifying application code.
  - ○ **Technical Dissection:**
    - ■ **mTLS (Mutual TLS):** Automatically encrypts and authenticates all traffic between Script Kitty's microservices, ensuring that only trusted services

can communicate. This is a crucial defense against eavesdropping and impersonation.

- **Fine-grained Access Control:** Policy enforcement at the network layer (e.g., "only Nexus can call Core's `process_query` endpoint, and only if authenticated"). This goes beyond simple network policies by understanding service identity.
- **Authorization Policies:** Define granular `AuthorizationPolicy` rules (e.g., `allow` or `deny` specific HTTP methods, paths, or source services).
- **Traffic Mirroring/Shadowing:** Allows for mirroring production traffic to a new version of a service for testing purposes without impacting live users, valuable for safely testing new agent behaviors.

- **6.4. Runtime Security (Falco):**

  - **Purpose:** Detects anomalous and suspicious behavior at the kernel level within Script Kitty's running containers, providing real-time threat detection.
  - **Technical Dissection:**
    - **System Call Monitoring:** Falco monitors system calls (`syscalls`) directly from the Linux kernel, providing deep visibility into container activity.
    - **Rule Engine:** Uses a declarative rule language to define suspicious behaviors (e.g., a process attempting to write to `/etc/passwd`, a container spawning an unexpected shell, a Skills agent trying to access unauthorized files, network connections to unknown IPs).
    - **Real-time Alerts:** Generates alerts (e.g., to Guardian, to a SIEM system, or to a Slack channel) when a rule is violated, enabling immediate response to potential compromises.
    - **Use Cases for Script Kitty:** Detects sandboxed escape attempts by Skills agents, unauthorized data exfiltration by Armory, or suspicious administrative activity within the Core or Foundry components.

- **6.5. Container Scanning (Trivy/Clair):**

  - **Purpose:** Scans Docker images for known vulnerabilities and misconfigurations before deployment.
  - **Technical Dissection:**
    - **Layer-by-Layer Scanning:** Analyzes each layer of a Docker image for vulnerabilities in operating system packages, language-specific dependencies (Python, Go, Node.js), and application libraries.
    - **Vulnerability Databases:** Integrates with public CVE databases (NVD) and commercial vulnerability intelligence feeds.
    - **SBOM (Software Bill of Materials) Generation:** Can generate a list of all components and dependencies within an image, improving transparency and auditability.

- - **Integration with CI/CD:** Scans are integrated into CI pipelines. Images with high-severity vulnerabilities can automatically fail the build, preventing their deployment.
    - **Use Cases for Script Kitty:** Ensures all base images and application dependencies for Nexus, Core, Armory, Skills, and Guardian are regularly scanned and patched.
- **6.6. Identity & Access Management (Keycloak):**

  - **Purpose:** Provides centralized authentication and authorization services for human users interacting with Script Kitty's administrative UIs (e.g., Guardian's oversight dashboard, Foundry's MLOps dashboards) and for internal service-to-service authentication.
  - **Technical Dissection:**
    - **OpenID Connect/OAuth2:** Implements standard protocols for secure authentication and authorization.
    - **User/Role Management:** Manages users, groups, and roles, defining what actions users or internal services are permitted to perform.
    - **Single Sign-On (SSO):** Enables human operators to log in once and access multiple Script Kitty administrative applications securely.
    - **Service Account Management:** Provides a robust way to manage credentials and access policies for internal Script Kitty service accounts (e.g., Foundry's CI/CD pipelines accessing Kubernetes API, Guardian's feedback system updating the GKGS).
    - **Integration with Kubernetes:** Service accounts are typically linked to Keycloak roles, allowing Kubernetes to automatically authorize internal pod communication based on their assigned identity.

**Foundry's Role in Script Kitty's Overall Functionality:**

- **"Seamless Transitions":** Achieved through robust CI/CD pipelines, GitOps-driven deployments (Blue/Green, Canary), and Kubernetes's rolling update capabilities, ensuring new models and features are deployed without disruption. Observability provides the confidence to make these transitions.
- **"Evolving":** Foundry is the enabler of continuous learning and improvement. It provides the data versioning, ML pipeline orchestration, and model serving infrastructure that Guardian uses to trigger and manage retraining cycles, and that Core uses to deploy new planning models.
- **"Universal Agent Features":** Every Script Kitty agent (Nexus, Core, Armory, Skills, Guardian) inherently benefits from Foundry's services: containerization, secure communication (mTLS via service mesh), centralized logging, metrics collection, secrets management, and data access. Foundry provides the standardized "platform" upon which all agents are built.
- **"Platform that evolves with all agents' knowledge":** Foundry provides the data infrastructure (Data Lake, DVC, Feature Store) for accumulating the collective experience and knowledge (e.g., Core's successful plans, Armory's extracted data,

Skills' task outcomes, Guardian's feedback). This accumulated data then feeds back into retraining cycles orchestrated by Foundry's ML pipelines, leading to the evolution of models across the entire Script Kitty system. The MLOps practices ensure that this knowledge is consistently applied and measured.

---

## Conceptual Code File Structure and Summaries:

This section outlines a high-level, conceptual file tree and provides a summary for key code files within each of Script Kitty's architectural components. This demonstrates the practical, microscopic layout of the codebase, aligning with the detailed technical analysis.

```
script-kitty/
├── README.md                    # Project overview and quick start guide.
├── CONTRIBUTING.md              # Guidelines for community contributions.
├── docs/                        # Comprehensive documentation, API references, architecture diagrams.
│   ├── architecture/
│   │   ├── overall_architecture.drawio
│   │   └── component_interactions.md
│   ├── apis/
│   │   └── grpc_api_reference.md
│   └── operational_guides/
│       └── deployment_guide.md
│
├── common/                      # Shared libraries, schemas, and utilities.
│   ├── proto/                   # Protobuf definitions for inter-service communication.
│   │   ├── script_kitty_messages.proto   # Core messages (user input, task requests/results).
│   │   ├── nexus_schemas.proto           # Nexus-specific schemas (session context, intent).
│   │   ├── core_schemas.proto            # Core-specific schemas (plans, agent assignments).
│   │   ├── armory_schemas.proto          # Armory-specific schemas (tool manifests, resource requests).
│   │   ├── skills_schemas.proto          # Skills-specific schemas (execution requests, results).
│   │   └── guardian_schemas.proto        # Guardian-specific schemas (policy proposals, metrics).
│   ├── python/
│   │   ├── __init__.py
│   │   ├── sdk_client.py        # Python gRPC client for internal communication.
│   │   ├── utils.py             # Common utility functions (logging, error handling).
│   │   └── exceptions.py        # Custom exception classes.
│   └── go/                      # Go-specific shared components (e.g., lightweight clients, helpers).
│       ├── sdk_client.go
│       └── utils.go
│
```

```
│
├── services/                    # Core microservices for each Script Kitty component.
│
│   ├── nexus-service/
│   │   ├── Dockerfile
│   │   ├── requirements.txt
│   │   ├── main.py              # Nexus microservice entry point, gRPC server setup.
│   │   ├── handlers/
│   │   │   ├── chat_handler.py      # Routes general chat, calls LLM for responses.
│   │   │   └── intent_router.py     # Logic for LLM-based and classifier-based intent routing.
│   │   ├── adapters/
│   │   │   ├── slack_adapter.py      # Handles Slack API integration (webhooks, message
sending).
│   │   │   ├── web_chat_adapter.py    # Manages WebSocket connections for web UI.
│   │   │   └── api_adapter.py        # Generic REST API for direct integrations.
│   │   ├── nlu/
│   │   │   ├── nexus_llm_inference.py  # Wrapper for vLLM/Ollama inference endpoint.
│   │   │   └── traditional_classifier.py # Scikit-learn/Transformers model for fast NLU fallback.
│   │   ├── context/
│   │   │   ├── session_manager.py     # Manages Redis/PostgreSQL session state.
│   │   │   └── conversation_history.py # Manages sliding window for LLM context.
│   │   └── preprocessors/
│   │       └── input_normalizer.py     # Unicode, whitespace, basic PII/toxicity filter.
│
│   ├── core-service/
│   │   ├── Dockerfile
│   │   ├── requirements.txt
│   │   ├── main.py              # Core microservice entry point, gRPC server setup.
│   │   ├── planning/
│   │   │   ├── htn_planner.py        # HTN decomposition logic, LLM prompting for plans.
│   │   │   └── plan_templates.py      # Loads and manages pre-defined plan templates.
│   │   ├── orchestration/
│   │   │   ├── agent_manager.py       # Interacts with Kubernetes API for agent lifecycle.
│   │   │   └── task_monitor.py        # Tracks progress of assigned tasks from other agents.
│   │   ├── knowledge_graph/
│   │   │   ├── gkgs_client.py        # Client for Neo4j/ArangoDB and Weaviate (vector store).
│   │   │   └── knowledge_updater.py   # Logic for updating GKGS based on feedback.
│   │   ├── communication/
│   │   │   ├── internal_message_bus.py # Kafka/Pulsar client for async communication.
│   │   │   └── grpc_server.py        # gRPC server implementation for Core's endpoints.
│   │   └── learning/
│   │       └── policy_refinement.py   # Placeholder for RL components for planning policy.
│
│   ├── armory-service/
```

```
│  │  │  ├── Dockerfile
│  │  │  ├── requirements.txt
│  │  │  ├── main.py                 # Armory microservice entry point.
│  │  │  ├── web_research/
│  │  │  │   ├── scraper_engine.py      # Playwright/Selenium integration for web navigation.
│  │  │  │   ├── data_extractor.py      # BeautifulSoup4, Trafilatura for HTML parsing/text
│  │  │  │                              extraction.
│  │  │  │   └── content_digester.py    # LLM summarization, NER/RE using
│  │  │  │                              spaCy/Transformers.
│  │  │  ├── tool_management/
│  │  │  │   ├── tool_registry_db.py    # Database interface for tool metadata (PostgreSQL).
│  │  │  │   └── tool_manifest_parser.py # Logic for ingesting tool documentation.
│  │  │  ├── wrapper_generation/
│  │  │  │   ├── code_llm_interface.py  # Interface to code-focused LLM for wrapper generation.
│  │  │  │   ├── template_renderer.py   # Jinja2 for injecting LLM code into secure templates.
│  │  │  │   └── human_review_api.py    # REST API for human verification of generated code.
│  │  │  ├── resource_provisioning/
│  │  │  │   ├── terraform_runner.py    # Executes Terraform IaC plans.
│  │  │  │   ├── ansible_runner.py      # Executes Ansible playbooks for configuration.
│  │  │  │   ├── cloud_sdk_client.py    # Direct integration with cloud provider SDKs (boto3,
│  │  │  │                              google-cloud-python).
│  │  │  │   └── quantum_client.py      # Qiskit, Cirq interfaces for quantum resource access.
│  │  │
│  │  ├── skills-service/
│  │  │  ├── Dockerfile
│  │  │  ├── requirements.txt
│  │  │  ├── main.py                 # Skills microservice entry point.
│  │  │  ├── executors/
│  │  │  │   ├── code_sandbox_executor.py# Manages Docker containers, Firejail, seccomp for
│  │  │  │                              code execution.
│  │  │  │   ├── cli_executor.py        # Securely runs CLI commands within sandboxes.
│  │  │  │   └── llm_code_generator.py  # CodeLlama/StarCoder2 interface for programming
│  │  │  │                              tasks.
│  │  │  ├── models/                 # Placeholder for specific ML model endpoints.
│  │  │  │   ├── computer_vision_model.py # TorchServe/TensorFlow Serving integration.
│  │  │  │   └── time_series_model.py
│  │  │  ├── data_analysis/
│  │  │  │   ├── pandas_processor.py    # Data manipulation with Pandas/Polars.
│  │  │  │   └── statistical_analyzer.py # SciPy, custom statistical routines.
│  │  │  ├── specialized_solvers/    # Modules for specific problem domains.
│  │  │  │   └── optimization_solver.py
│  │  │  ├── config/
│  │  │  │   └── sandbox_security_profiles.yaml # Sandbox resource limits and security rules.
│  │  │
```

```
├── guardian-service/
│   ├── Dockerfile
│   ├── requirements.txt
│   ├── main.py                    # Guardian microservice entry point.
│   ├── policy_enforcement/
│   │   ├── opa_client.py          # Client for Open Policy Agent.
│   │   └── custom_rule_engine.py  # (Optional) More complex ethical rule processing.
│   │   └── ethical_llm_evaluator.py# Fine-tuned LLM for ethical flag detection.
│   ├── monitoring/
│   │   ├── metrics_consumer.py     # Consumes Prometheus metrics streams.
│   │   ├── log_analyzer.py         # Processes logs from Loki/Elasticsearch for anomalies.
│   │   └── anomaly_detector.py     # Implements PyOD for runtime anomaly detection.
│   ├── evaluation_learning/
│   │   ├── feedback_processor.py   # Ingests human feedback, stores in data lake.
│   │   ├── evaluation_pipeline_trigger.py # Triggers Kubeflow/Argo Workflows for retraining.
│   │   └── bias_mitigation.py      # Integrates Aequitas, Fairlearn, AIF360 for bias analysis.
│   ├── user_interface/
│   │   └── dashboard_api.py        # Provides data for Streamlit/Plotly Dash UI for human
oversight.
│   ├── policies/
│   │   ├── security_policies.rego  # Rego policies for OPA enforcement.
│   │   └── ethical_guidelines.md   # Human-readable ethical principles.
│
├── infrastructure/                # IaC and deployment manifests.
│   ├── kubernetes/
│   │   ├── base/
│   │   │   ├── namespace.yaml
│   │   │   ├── common_services.yaml    # Prometheus, Grafana, Loki, Vault, Keycloak
deployments.
│   │   │   ├── ingress_controller.yaml
│   │   │   └── network_policies.yaml    # Calico/Cilium base policies.
│   │   ├── overlays/                # Kustomize/Helm values for different environments (dev,
staging, prod).
│   │   │   ├── dev/
│   │   │   │   ├── kustomization.yaml
│   │   │   │   └── resource_limits.yaml.patch
│   │   │   └── prod/
│   │   │       ├── kustomization.yaml
│   │   │       └── hpa_configs.yaml
│   │   ├── crds/                    # Custom Resource Definitions for Script Kitty agents.
│   │   │   ├── scriptkitty_agent.yaml
│   │   │   └── skill_module.yaml
│   │   └── services/                # Specific Kubernetes manifests for each microservice.
│   │       ├── nexus/
```

```
│   │   │       ├── deployment.yaml
│   │   │       ├── service.yaml
│   │   │       └── hpa.yaml
│   │   │   ├── core/
│   │   │   │   ├── deployment.yaml
│   │   │   │   └── service.yaml
│   │   │   └── # ... and so on for armory, skills, guardian ...
│   │   ├── terraform/              # Terraform modules for cloud infrastructure provisioning.
│   │   │   ├── aws/
│   │   │   │   ├── main.tf          # Main Terraform configuration for AWS.
│   │   │   │   ├── vpc.tf           # VPC network configuration.
│   │   │   │   └── eks_cluster.tf   # EKS (Kubernetes) cluster provisioning.
│   │   │   └── gcp/
│   │   │       ├── main.tf
│   │   │       └── gke_cluster.tf   # GKE (Kubernetes) cluster provisioning.
│   │   ├── ansible/                # Ansible playbooks for machine configuration.
│   │   │   ├── playbooks/
│   │   │   │   └── node_setup.yaml      # Configuring worker nodes (e.g., GPU drivers).
│   │   │   └── roles/
│   │   │       └── security_hardening/
│   │   │           └── tasks/main.yaml
│   │
│   ├── pipelines/                  # CI/CD and MLOps pipeline definitions.
│   │   ├── gitlab-ci/              # GitLab CI/CD pipeline configurations.
│   │   │   └── .gitlab-ci.yml
│   │   ├── github-actions/         # GitHub Actions workflow definitions.
│   │   │   └── .github/workflows/main.yaml
│   │   ├── kubeflow/               # Kubeflow Pipelines definitions for ML workflows.
│   │   │   ├── training_pipeline.py     # Python DSL for defining model training workflow.
│   │   │   ├── evaluation_pipeline.py
│   │   │   └── model_deployment_pipeline.py
│   │   ├── argo_workflows/         # Argo Workflows definitions.
│   │   │   ├── data_ingestion.yaml
│   │   │   └── model_retraining.yaml
│   │   ├── dvc/                    # DVC pipeline definitions.
│   │   │   └── dvc.yaml            # Defines data processing stages (data, featurization, model).
│   │   └── airflow/                # Apache Airflow DAGs for data orchestration.
│   │       └── dags/
│   │           └── web_data_etl.py      # DAG for scraping and processing web data.
│   │
│   ├── data/                       # Data lake structure (S3/GCS buckets, DVC tracking).
│   │   ├── raw/                    # Raw, immutable ingested data.
│   │   │   ├── user_inputs/
│   │   │   └── web_scraping/
```

```
│   ├── processed/              # Cleaned, pre-processed data.
│   │   └── nlu_training_data/
│   ├── features/               # Features for ML models (tracked by Feast/DVC).
│   │   └── user_embeddings/
│   └── models/                 # Trained model artifacts (tracked by DVC).
│       ├── nexus_llm/
│       ├── core_planner/
│       └── skills_vision_model/
│
└── tools/                      # Development and operational scripts.
    ├── dev_setup.sh            # Script for setting up a local development environment.
    ├── deploy_local.sh         # Script for deploying a local Kubernetes instance
(Minikube/Kind).
    ├── cli/                    # CLI tools for interacting with Script Kitty.
    │   └── sk_cli.py           # Python CLI for admin tasks, e.g., `sk_cli core deploy-agent`.
    └── scripts/
        ├── backup_db.sh
        └── monitor_alerts.py
```