

System Architecture Document – Elspeth

Executive Summary

Elspeth is a Python-based framework for orchestrating Large Language Model (LLM) experiments in a secure and extensible manner ¹. The system is designed as a **command-line tool** that loads input datasets, runs multiple LLM prompt trials (experiments), and produces results and reports under strict compliance controls. **Key architectural features include** a plugin-driven design (for datasources, LLM integrations, metrics, and outputs), an emphasis on **security-by-design** (classified data handling, prompt sanitization, cryptographic signing of outputs), and a flexible configuration system for defining experiment suites.

The architecture comprises **distinct layers**: a CLI interface for user input, a core orchestration engine (`ExperimentOrchestrator`) that bridges data sources, LLM clients, and result sinks ², and a set of plugin interfaces that encapsulate domain-specific logic (e.g. data loading, LLM API calls, post-processing metrics). Data flows from configured **datasources** (e.g. CSV files or Azure Blob Storage) through the orchestrator into LLM prompts and then to various **output sinks** (CSV/Excel files, repository commits, etc.), with optional **analysis plugins** computing metrics and comparisons along the way. Multiple experiments can be grouped into suites for back-to-back execution and baseline comparisons via the `ExperimentSuiteRunner` component.

Through code analysis, this document confirms that Elspeth's implementation aligns with its stated goal of **responsible, auditable LLM experimentation** ¹. Each input record is processed in isolation (or batched with controlled concurrency) to generate LLM outputs, which are logged and aggregated with rich metadata. Notably, the system enforces **data security levels** on all artifacts and output channels: every plugin and sink is tagged with a security classification (e.g. "OFFICIAL") and the framework prevents high-classification data from flowing into lower-cleared outputs ³. This ensures compliance with data handling policies. Additionally, all prompt content rendering uses a strict Jinja2 engine to avoid template injection or accidental omission of required fields ⁴ ⁵, and output files (like spreadsheets) are sanitized to mitigate injection attacks (e.g. Excel formulas are neutralized by the CSV sink) as part of the system's "security by default" posture ⁶.

Performance & scalability: The codebase supports concurrent processing of inputs via multithreading, gating parallelism based on workload size and external rate limits ⁷ ⁸. A rate-limiter plugin can throttle LLM API calls to respect service quotas ⁹, and a cost-tracker can accumulate token or API usage cost metrics across runs. These controls enable resilient execution under production-like conditions ¹⁰. While horizontal scaling (across machines) is not built-in, the modular architecture could be deployed in distributed settings or containerized for larger workloads. Reliability features include **retry logic** for LLM calls with exponential backoff ¹¹ ¹² and an optional checkpointing mechanism to resume long runs if interrupted ¹³.

Maintainability & extensibility: Elspeth's architecture cleanly separates concerns and uses abstract protocols for key components (data sources, LLM clients, result sinks, etc.), making it straightforward to add new plugins without altering core logic. The plugin registry allows registering custom plugins at runtime ¹⁴ ¹⁵, reflecting a conscious design decision to favor composition over monolithic design. The codebase is thoroughly unit-tested (with over 50 test modules covering everything from CLI to

plugins ¹⁶ ¹⁷), which bolsters confidence in changes and aids accreditation. Documentation (though partially outdated in `/docs`) and in-code dataclass definitions provide clarity on configuration schema, and an internal traceability matrix maps security controls in code to compliance requirements.

Key Findings and Recommendations: The analysis confirms that Elspeth implements robust security and compliance measures (classification enforcement, audit logging hooks, output signing) suitable for an Authority To Operate (ATO) in regulated environments. Some areas for improvement were identified: for instance, exception handling around certain plugin operations (e.g. aggregator finalize steps) could be hardened to avoid run interruption on plugin errors (see Risk Register). Additionally, the project currently lacks a finalized open-source license and some documentation lags behind recent refactors, which should be addressed to facilitate secure adoption and oversight. Section 9 of this document provides a detailed list of critical issues and recommended mitigations before deployment in a production or accredited environment.

1. System Context and Scope

1.1 System Purpose

Elspeth is intended to **enable secure, controlled experimentation with LLMs** – allowing users (typically ML engineers or researchers) to run comparative studies of prompts and models while ensuring the process is auditable and compliant ¹. By examining the source code, we infer the system’s purpose is to streamline the execution of experiments (comprising one or more LLM prompt evaluations) and produce structured outputs (metrics, reports) without exposing sensitive data or violating usage policies. Key capabilities include: - **Flexible experiment definition:** Users can define prompts, variations (“criteria”), and metrics via configuration, then execute them at scale through the CLI ¹⁸ ¹⁹. - **Security and compliance:** Every step from input ingestion to result output is designed with security controls (e.g. prompt sanitization, mandatory classification labels, cryptographic signing) to support compliance requirements ⁶ ³. - **Modularity:** The system supports “drop-in” components (plugins) for data sources, LLM backends, processing middleware, and output sinks ²⁰, reflecting a purpose of serving as a **framework** rather than a single-purpose script. For example, organizations can plug in their own data connector or a different LLM API without modifying the orchestrator logic ²⁰. - **Reproducibility and audit:** The code places emphasis on determinism and traceability. Prompts are rendered with a strict template engine to avoid nondeterministic formatting ⁴ ⁵, and the execution pipeline can record each attempt and outcome (including retries, timestamps, and metrics) for later review ²¹ ²². Optional **signing of result bundles** further ensures integrity for compliance audits (so outputs can be verified as unaltered).

In summary, the system’s purpose is to let teams **experiment with LLMs “without compromising compliance or auditability”** ¹. Code evidence of this mission includes the built-in **policy enforcement hooks** (e.g., `security_level` tags on data and sinks ³) and the careful handling of content (e.g., automatic formula sanitization in output spreadsheets to prevent injection attacks ²³).

1.2 System Boundaries

Elspeth operates primarily as a **command-line application**, which defines its system boundary at the process level. The **entry point** is the CLI module (`elspeth.cli`), which parses user arguments and invokes the orchestrator logic ²⁴ ²⁵. There is no continuously running server or service in the architecture – each invocation of the CLI is a self-contained execution of one or more experiments. The absence of any listening network sockets or web handlers in the code confirms that **Elspeth does not**

expose a network API or interactive UI; instead, it reads configuration files and data, then produces outputs to files or external systems, all within the lifecycle of a CLI process.

Inputs (system ingress):

- **Configuration files:** Typically a YAML file (e.g. `settings.yaml`) provides orchestrator settings, including data source info, LLM credentials or endpoints, and plugin configurations. The CLI accepts a `--settings` argument to specify this file ²⁴. The code loads and validates this config into a `Settings` dataclass instance, ensuring required sections are present ²⁶.
- **Experiment data:** The primary data input is usually tabular data (prompts and any associated fields) either from a CSV file or an external source (like Azure Blob storage). This is abstracted by the `DataSource` interface. For example, a `CSVDataSource` plugin will read a local CSV, whereas a `BlobDataSource` will fetch data from cloud storage. The orchestrator then invokes `datasource.load()` to ingest this data as a pandas `DataFrame` ²⁷ ²⁸.
- **User command arguments:** The CLI arguments define the execution mode (single experiment vs suite, profiles, etc.). For instance, `--single-run` toggles between running one experiment versus a suite ²⁹, and `--validate-schemas` triggers a dry-run mode for schema checking only ³⁰ ³¹. These inputs affect control flow but are all handled at process start (no interactive runtime input beyond this). Authentication credentials for external services may also enter as environment variables (for Azure or OpenAI keys), but those are consumed by underlying SDKs (`azure-identity`, `openai`) rather than direct CLI input.

Outputs (system egress):

- **Result files:** By default, results are written to designated output files or directories. For example, the CLI can output the combined result dataset as a CSV (`--output-csv`) ³², and write analytics reports (CSV, Excel, charts) to an output directory (`--reports-dir`) ³³. Individual **ResultSink** plugins handle their own output logic: e.g., `CsvResultSink` writes rows to CSV files, `ExcelSink` produces an Excel workbook, `VisualReportSink` generates images/HTML charts, etc. All sinks are invoked at the end of a run via the artifact pipeline (described later), and they operate either on the local filesystem or through external APIs (e.g., a GitHub repository sink uses network calls to create a commit).
- **External service calls:** Some outputs or intermediate actions cross the system boundary via APIs. For instance, a **repository sink** may call GitHub or Azure DevOps REST APIs (over HTTPS) to upload results to a remote repository ³⁴. Similarly, if using Azure Blob as an output sink (if configured), the system will use Azure SDK calls to write files to cloud storage. These external calls are considered part of output handling.
- **Logs and console:** Operationally, Elspeth emits logging info to the console (or stdout). For example, it logs summary info for each experiment's row count ³⁵ and any warnings from validation. These logs can be captured by the environment or redirected but are not stored by Elspeth beyond in-memory during execution. The log messages (via Python's logging) provide an audit trail of key events (start/stop, warnings, errors) ³⁶ ³⁷.

External system interactions fall outside Elspeth's trust boundary and are treated via plugins:

- Calls to LLM services (OpenAI/Azure OpenAI) are made via the LLM client plugins and the `requests` or official SDK libraries ³⁸. These are outbound HTTPS calls and require API credentials, but the details (URL, keys) are configured in the YAML or environment.
- Data fetching from cloud (Azure Blob) similarly goes through Azure SDK; the code uses `azure-identity` to handle auth and `azure-storage-blob` to stream data ³⁹.
- No incoming network connections: Elspeth does not listen for requests or events. Any **automation** would be via external schedulers or CI pipelines invoking the CLI with appropriate parameters.

In summary, the system boundary encloses the **Elspeth CLI process** and its direct file/system I/O. The **entry point** is the CLI's `main()` function, which processes arguments and calls `ExperimentOrchestrator.run()` for single experiments ²⁵ or `ExperimentSuiteRunner.run()` for suites ⁴⁰ ⁴¹. The boundary with the user is at the CLI (the

user or a script triggers the run and provides config), and the boundary with external services is at the plugin integration points where HTTP SDK calls are made. All such calls are outbound; there are no backcalls or webhooks coming into Elspeth.

1.3 Stakeholders and Actors

Because Elspeth is not a multi-user server application, “actors” in the traditional sense (end-users with roles interacting concurrently) are limited. However, we can identify the **primary user role** and **external systems** as actors in the context of the system:

- **Experiment Operator (Primary User):** This is typically a data scientist, ML engineer, or developer who uses Elspeth to run experiments. They interact with the system by preparing configuration files and data, then executing the CLI. In code, this actor’s intent is represented by the CLI arguments and configuration. The operator can choose to run a single experiment or a suite, enable or disable certain plugins (e.g. `--disable-metrics` flag to turn off metrics plugins⁴²), and direct outputs to certain locations. This actor trusts the system to enforce the rules (like not leaking data) and to produce the promised reports. There is no authentication or authorization mechanism *within* Elspeth for multiple user accounts – it assumes the user running the CLI is authorized to access the provided data and configurations (any user-level access control would be handled by the OS or environment outside the application).
- **System Integrator/Administrator:** Although not explicitly coded, a stakeholder exists in setting up Elspeth in an environment (e.g., adding custom plugins, ensuring compliance settings). This might be a platform engineer who configures certain defaults or integrates Elspeth into pipelines. For example, they might pre-register internally developed plugins via the `plugin_registry` before execution, or set environment variables for API keys. This stakeholder ensures that Elspeth runs with the appropriate context (Azure credentials in place, etc.).
- **External LLM Service** (e.g., OpenAI API or Azure OpenAI Service): These external systems act as *actors* from Elspeth’s perspective in that the orchestrator will send requests to them and receive responses. The LLM client plugins encapsulate this interaction. For instance, the `OpenAIClient` plugin uses the OpenAI SDK to send the prompt and get a completion³⁸. The availability and performance of these external actors directly affect Elspeth’s operation. Elspeth’s design acknowledges them via retry logic and rate limiting around calls^{11 12}. These services require credentials (API keys or Azure identity), which the operator or admin must supply via config or environment. There is an implicit trust that the external LLM service will handle the prompts per its specification (with any content filtering as needed). If the external service fails or returns an error, Elspeth will capture that as a failure record and optionally retry^{43 44}.
- **Data Storage Services:** If using an external data source or sink (e.g., Azure Blob Storage for input, or an enterprise Git repository for outputs), those act as external actors too. The `BlobDataSource` plugin, for example, uses Azure’s blob client to download a file; the blob storage is the actor providing data. Similarly, a `GitHubRepoSink` acts by pushing data to a GitHub repository (the GitHub service is the actor receiving data). These interactions are authenticated via tokens or identities and occur over secure channels. Elspeth treats them as plugins – for instance, when a repository sink is configured, the system will prepare an HTTP payload to the GitHub API in the `outputs/repository.py` logic (not shown here, but the use of `requests` for repository sinks is noted⁴⁵). The stakeholder interest here is that these systems receive only the data they should (hence the security level checks before sending, see

Section 4.2 Security) and that failures in these systems don't break the orchestrator (Elspeth marks output failures but continues execution of other outputs).

- **Compliance/Audit Team:** While not directly interacting with Elspeth at runtime, this stakeholder consumes the outputs and logs. They rely on Elspeth's features like signing of results and detailed metadata to audit what prompts were run, what responses were received, and whether any policy rules were violated. For example, if required by regulation, a compliance officer might use the `TRACEABILITY_MATRIX.md` and the logs/artifacts to verify that every piece of content is tagged with a classification and that outputs have not been tampered (signature verification). This stakeholder influences requirements (leading to features like `security_level` and `determinism_level` tracking in code ⁴⁶ ³), even though they are not driving the CLI.

In terms of **roles or permissions** within Elspeth: because the execution is single-user, the code doesn't implement role-based access control. However, it does implement *data access control* via the aforementioned security levels. You can think of each plugin/sink as having a role of handling data up to a certain sensitivity. The function `is_security_level_allowed(producer_level, consumer_level)` explicitly enforces a simple policy: e.g., a sink with clearance "Official" cannot receive "Secret" data from a producer ³. This acts as an **internal authorization check** for data flows rather than user actions. In tests, for instance, sinks are tagged `_elspeth_security_level = "official"` ⁴⁷, and the orchestrator computes an experiment's overall security classification from its inputs ⁴⁸. These levels correspond to organizational data policies, so one might say the "Security Officer" is an implicit stakeholder who sets what those levels mean (likely via config) and ensures Elspeth's enforcement meets the required standard.

1.4 External Dependencies and Integrations

Elspeth integrates with several external systems and libraries, which we identify from the code and configuration:

- **Python Runtime and Libraries:** As a Python application, Elspeth depends on numerous libraries for core functionality. Notably, it uses **pandas** for data handling (DataFrame operations) ⁴⁹, **PyYAML** for reading YAML configs ⁴⁹, **Jinja2** for prompt templating ⁵⁰, and **jsonschema** for validating configuration structures ⁵¹. These are internal dependencies but critical for operations (e.g., the `PromptEngine` wraps Jinja2 to render prompt templates ⁴, and config validation likely uses jsonschema under the hood for schema definitions). In addition, scientific libraries like **scipy** and **statsmodels** are optional dependencies for advanced statistical analysis plugins ⁵².
- **LLM Service APIs:** Elspeth can call out to external LLM services. The **OpenAI API** is one such integration – the project includes the `openai` Python library (OpenAI SDK) as a dependency ⁴⁹. In the code, there's an `openai_http` plugin (likely sending requests via the OpenAI library or `requests` if needed). For Azure's variant of OpenAI (Azure OpenAI Service), Elspeth uses Azure's identity and likely calls OpenAI endpoints with an API version. The `azure_openai.py` plugin and `azureml-core` (for Azure ML telemetry middleware) indicate integration with Azure's ecosystem ⁵³ ⁵⁴. For example, azure-identity is used to obtain credentials seamlessly (Managed Identity or Service Principal) ⁵⁵, which the Azure OpenAI plugin can utilize instead of static API keys. All interactions with these LLM services are **outbound HTTPS calls** (through the libraries), and Elspeth does not directly expose these keys or endpoints beyond passing them to the SDK. The code wraps these calls in middleware and rate-limiters: e.g., before calling `llm_client.generate()`, it may inject a delay via `rate_limiter.acquire()` ⁹ and

after calling, it passes the response through any `LLMMiddleware.after_response` hooks (for logging or redaction) ⁵⁶ ⁵⁷ .

- **Cloud Storage:** If configured to use Azure Blob Storage as a data source or output sink, Elspeth relies on **Azure SDKs**:

- The `azure-storage-blob` library is listed as a dependency ⁴⁹ . Code in `datasources/blob_store.py` (as referenced in docs) uses this to fetch data from containers ³⁹ . The CLI's docstring explicitly mentions hydrating experiment data from Azure Blob with `config/blob_store.yaml` ⁵⁸ . This implies the config can specify a blob container and file, and Elspeth will use `azure-storage-blob` to download it. The `azure-identity` integration means credentials are not hard-coded but pulled from the environment or managed identity ⁵⁵ .
- As an output, if a sink plugin writes to Azure Blob, it would also use `azure-storage-blob` to upload. The code ensures any such sink runs only if TLS is in place (implicitly, Azure SDK uses HTTPS by default; the docs note to ensure certificate verification is handled ⁵⁹).

- These interactions with Azure happen behind the scenes in plugins; from Elspeth's perspective, they are atomic operations like any file read/write, except they can fail due to network issues or auth. The system should log those failures if they occur (likely as exceptions captured in plugin errors).

- **GitHub or DevOps Repositories:** A unique integration is the **repository sinks** (plugins named `"github_repo"` or `"azure_devops_repo"`). The CLI adjusts a `dry_run` flag for these if live outputs are disabled ³⁴ . That suggests that by default, repository sinks might operate in a "dry run" mode (just simulate or log what they *would* do) unless `--live-outputs` is passed. When live, they will use API calls:

- The GitHub repository sink presumably uses GitHub's REST API (with `requests`) to create a commit or upload a file. Indeed, the dependency analysis notes that `requests>=2.31.0` is used by *repository sinks* among other things ⁴⁵ . The code likely constructs an HTTP PUT/POST to the GitHub content API with a token provided in config.
- The Azure DevOps repository sink similarly would call Azure DevOps REST endpoints for pushing files. Both require tokens or PATs configured in the sink options.
- In terms of system boundaries, these are outbound calls. They raise an interesting security point: pushing to a repository could be equivalent to publishing potentially sensitive data. That's why the code explicitly includes these in the dry-run control – by default, they won't actually push unless the user allows live outputs ³⁴ . Also, classification enforcement is crucial here (the pipeline will prevent, say, "SECRET" data being pushed to a repo sink not cleared for it ³).

- **Email/Notification:** There is no evidence in the code of email or messaging integrations. No SMTP or similar library is present in dependencies. So notifications are not directly sent by Elspeth except via the logs and outputs.

- **Visualization and Analysis:** For generating charts and statistical analyses as output artifacts, Elspeth uses libraries like **matplotlib** and **seaborn** (for visualizations) and **pingouin/statsmodels** (for statistical tests) as optional extras ⁶⁰ ⁵² . These aren't external services, but they are noteworthy external libraries that produce files or figures included in reports. For example, the `VisualReportSink` likely uses matplotlib to generate a PNG or HTML with embedded images ⁶¹ . These images might be later opened in a browser or office tool by users, which means any data in them should be properly sanitized to avoid XSS or similar (for instance,

if a chart had textual annotations from input data, it shouldn't include malicious HTML – although using matplotlib means output is an image, which is inherently safe). The code sets up these libs only when needed (they are extras), showing consideration of minimizing attack surface by not installing them unless required ⁶².

- **Operating Environment:** While not a dependency per se, the environment (OS, Python version) matters. The README calls for Python 3.12 and a Unix-like shell for Makefile usage ⁶³. There's also mention of `pytype` and `ruff` for linting in dev dependencies ⁶⁴, but those don't affect runtime. The system likely assumes a standard file system and possibly internet connectivity to reach external APIs. In high-security deployments, internet may be restricted; Elspeth can still function with the **mock LLM** or offline mode for testing (the `mock.py` plugin provides a local LLM stub ⁶⁵). The existence of a mock LLM plugin means one can run experiments without any external calls (useful for offline or testing scenarios), which is an integration with a “dummy” actor that always returns deterministic content ⁶⁵.

In summary, **Elspeth's external integrations** include **cloud services (Azure, OpenAI)**, **version control systems (GitHub/Azure DevOps)**, and a range of third-party libraries for functionality. All these interactions are handled through well-defined plugin interfaces or libraries, keeping the core orchestrator agnostic of the specifics. This design ensures that, for example, switching from OpenAI to an internal model only requires writing a new LLM client plugin implementing `generate()`, or switching data storage from Azure Blob to, say, AWS S3 would involve a new `datasource` plugin. The plugin registries and factory functions (e.g., `llm_registry.create(name, options)`) serve as the abstraction layer to instantiate the correct integration at runtime ⁶⁶ ⁶⁷.

From a **security perspective**, external calls are made over TLS and rely on the security of external systems (we'll discuss later how credentials are managed and how data is classified to protect these calls). The code's use of updated library versions (requests 2.31.0, azure-storage 12.19.0, etc.) indicates an effort to stay current with security patches ⁴⁹. Notably, the dependency analysis file suggests monitoring CVEs for these integrations (e.g. Azure SDK credential issues, OpenAI library logging vulnerabilities) ⁶⁸ ⁴⁵, which is an operational task for whoever maintains the deployment.

2. Architectural Views

2.1 Logical Architecture (Layered View)

Elspeth's logical architecture can be described in layers and components that interact to achieve the overall functionality. By reading the code, we identify the following **primary layers**: - **Presentation Layer (CLI Interface)** – Handles user interaction, input parsing, and high-level output presentation on the console. - **Orchestration & Control Layer** – Coordinates the experiment flow, including reading data, invoking LLM calls, applying plugins, and triggering output generation. - **Plugins Layer (Extensible Components)** – Encapsulates domain-specific or customizable logic in several categories: Data Sources, LLM Clients & Middleware, Experiment Plugins (metrics, validators, early-stoppers), and Output Sinks. - **Data/Storage Layer** – Deals with data persistence and access: reading input data from files or databases, and writing output artifacts to files or external stores. (Note: Elspeth doesn't use a traditional database, but this layer covers file I/O and external storage via plugins.)

Each layer is largely separated by abstract interfaces (protocols), which can be seen in the code. For example, the **CLI (presentation)** uses only the exposed interfaces of the orchestrator and settings loader to do its job – it does not manipulate internal data structures directly but calls

`validate_settings()`, `load_settings()`, then `ExperimentOrchestrator.run()` or `ExperimentSuiteRunner.run()` ²⁶ ²⁵.

Layer breakdown with examples from code:

- **Presentation (CLI):** This is implemented in `elspeth.cli.py`. It defines an argument parser to accept user inputs like config file path, verbosity, toggles for metrics or outputs, etc. ²⁴ ⁶⁹. It prints previews of results for the user and handles writing a quick output CSV if requested ³² ⁷⁰. The CLI layer also logs warnings or errors in a user-friendly way (e.g., after running a suite, it logs each experiment's result count ³⁵ and outputs any validation warnings). Essentially, CLI is the thin wrapper that **interprets user intent and calls into the core layers**. There is no graphical UI or web UI – the CLI is the sole presentation interface.
- **Orchestration (Core):** This corresponds to classes in `elspeth.core` such as `ExperimentOrchestrator` and `ExperimentSuiteRunner`. The orchestrator is the central coordinator for a single experiment's execution: it ties together a DataSource, an LLM client, and a list of sinks, along with all plugin definitions from the config ²⁵. When `ExperimentOrchestrator.run()` is invoked, it loads the data and then delegates to an `ExperimentRunner` which handles the row-by-row prompt execution logic ⁷¹ ⁷². The orchestrator itself handles setup: for instance, it creates instances of all required plugins (row, aggregator, validation, etc.) by calling factory functions in the plugin registry ⁷³ ⁷⁴. It also sets up a **PluginContext** with a resolved security level and passes it down so that every plugin and sink knows the context of the experiment (name, provenance, and allowed security level) ⁴⁸ ⁷⁵. This layer also includes the **control components** such as `RateLimiter` and `CostTracker` – these are instantiated based on config (using factories `create_rate_limiter` and `create_cost_tracker`) and attached to the orchestrator/runner if present ⁷⁶ ⁷⁷. They function as cross-cutting concerns, limiting throughput and tracking usage globally across the run.

The `ExperimentSuiteRunner` extends the orchestration layer to handle multiple experiments. It loads the suite configuration (which may include a baseline and several experiments), merges defaults with each experiment's specific config via the `ConfigMerger`, and then essentially creates an `ExperimentRunner` for each experiment sequentially ⁷⁸ ⁷⁹. The suite runner is responsible for baseline comparisons – after running all experiments, it invokes any baseline comparison plugins to compare each experiment's results to the baseline's results ⁸⁰ ⁸¹. It uses plugin context similarly and reuses any shared LLM middleware between experiments for efficiency ⁸². This orchestration layer ensures all experiments in a suite are run in a controlled way, and it gathers their outputs into a combined results structure for reporting ⁸³.

- **Plugins (Extensible Logic):** This logical layer is subdivided into plugin types, each with its interface:
- **Data Source Plugins:** Provide the data to run experiments on. They implement a simple interface: a `load()` method returning a pandas DataFrame ⁸⁴. In config, a datasource is specified by name and options (e.g., `{"plugin": "csv_local", "options": {"path": "data.csv"}}`). The code maps this to a concrete class (e.g., `CSVLocalDataSource`) via a registry. Once instantiated, the orchestrator calls `datasource.load()` to get the DataFrame ⁸⁵. DataSource plugins also can tag the DataFrame with metadata like schema or security level. For example, the `BlobDataSource` might tag data with a certain security level if loaded from a secure container. The CLI focuses heavily on the data source – its docstring indicates initial emphasis on Azure Blob input ⁵⁸.

- **LLM Client Plugins:** These are how the system interfaces with language models. An LLM plugin implements `LLMClientProtocol`, essentially a `generate(system_prompt, user_prompt, metadata) -> response` method ⁸⁶. For example, `OpenAIClient` likely wraps `openai.Completion.create()`, and `AzureOpenAIClient` might call Azure's REST endpoint or the OpenAI SDK with an endpoint override. The code uses the term "LLM adapters" in docs ²⁰, and indeed in `plugins/llms` we have files like `openai_http.py`, `azure_openai.py`, `mock.py`. Each such plugin has a name (for config), and the `llm_registry.create(plugin_name, options)` returns an instance. The orchestrator receives this instance as `llm_client` and passes it into the runner. The **LLM middleware** are related plugins that wrap around the LLM client calls (see below).
- **Experiment Workflow Plugins:** These run during or after experiment execution to add custom logic:
 - **Row Plugins** operate on each row's result. They must implement a method (commonly `process_row(row, responses) -> metrics`). In code, after an LLM response is obtained for a row, the runner calls each row plugin to process that result ⁸⁷. For instance, a row plugin might extract a score from the LLM's content or perform a classification. The example in tests defines a row plugin that simply returns a constant metric ⁸⁸. These metrics are attached to the result record (e.g., `payload["results"][i]["metrics"]["custom_metric"] = 7` as seen in the test) ⁸⁹.
 - **Aggregator Plugins** operate after all rows are processed, to compute experiment-level aggregates. They implement e.g. `finalize(records) -> aggregate_metrics`. The runner iterates through aggregator plugins once all results are collected ⁹⁰. For example, a plugin might calculate the average score over all rows or count of a certain category. In the test, an aggregator plugin returns a count of records ⁹¹, which the runner stores under `payload["aggregates"][plugin_name]` ⁹². Those aggregates also propagate to sink metadata for reporting ⁹³.
 - **Validation Plugins** likely run on each LLM response to validate content or format. The runner calls `self._run_validations(response, request, row_context)` after getting a response ⁹⁴. Although the internal code isn't fully shown, we can infer that validation plugins might check if the response meets certain criteria (e.g., does JSON parse, does it contain prohibited content). If a validation fails, they might raise a `PromptValidationError` or similar, which the runner catches as a failure for that row ⁹⁵. One built-in validation appears to be **schema validation** – the system can attach a schema to the DataFrame (`df.attrs['schema']`) and check plugin compatibility ⁹⁶. This could ensure that if a plugin expects certain columns, the datasource provided them (failing early if not).
 - **Early-Stop Plugins** check conditions to decide whether to halt further processing. They implement a method (e.g., `check(record, metadata) -> reason or None`). The runner integrates this by evaluating after each successful row: it passes the latest record to `plugin.check()`, and if any returns a reason (like "cost_limit_exceeded" or "quality_threshold_met"), the runner sets an event to stop processing more rows ⁹⁷ ⁹⁸. This mechanism is akin to a circuit breaker – for instance, an early-stop plugin named "threshold" could stop the run if, say, 10 failures have occurred or if a certain accuracy metric is reached early. The code ensures thread-safe access by locking around early-stop checks, to avoid race conditions in parallel mode ⁹⁹. In config, early-stop conditions can be given as a simple dict (like `{"failure_rate": 0.5}`) which is normalized to a plugin definition ¹⁰⁰.
 - **Baseline Comparison Plugins** apply only in suite context, after baseline and another experiment have completed. They implement a `compare(baseline_payload, experiment_payload) -> diff` method. The suite runner calls them for each non-baseline experiment ¹⁰¹. This can produce comparative metrics (for example, "delta in

success rate: +5%"). The results are stored in the experiment's payload under `"baseline_comparison"` ¹⁰². Baseline plugins allow formal statistical comparison or differences (perhaps one plugin does a t-test on numeric scores).

- **LLM Middleware:** These act as interceptors around LLM calls. The orchestrator creates them via `create_middlewares` for any `llm_middleware_defs` in config ¹⁰³. Middleware can implement hooks such as `before_request(request)` to modify or log the request, `after_response(request, response)` to inspect or transform the response, and even suite-level hooks (`on_experiment_start`, `on_experiment_complete`, etc.) ¹⁰⁴ ¹⁰⁵. For instance, a middleware might implement OpenAI *content filtering* by sending the response to an Azure Content Safety API; or a telemetry middleware might log usage to Azure ML. The code explicitly calls these hooks at appropriate times: before each LLM call, each middleware's `before_request` is applied in order ¹⁰⁶, and after getting a response, `after_response` is applied in reverse order ⁵⁷ (to allow stacking effects). The suite runner also triggers `on_suite_loaded`, `on_experiment_start/complete`, etc., to let middleware gather info or adjust settings across experiments ¹⁰⁷ ¹⁰⁵. Middleware are a powerful way to incorporate enterprise concerns like logging, censoring, or dynamic prompt adjustments without changing the core LLM plugin logic.
- **Result Sink Plugins:** These handle output of results and artifacts. Each sink plugin implements the `ResultSink` protocol, primarily a `write(results, metadata)` method that the framework calls after obtaining the full experiment or suite results ¹⁰⁸. Examples include:
 - `CsvResultSink`: writes the list of result records to a CSV file. Indeed, in CLI when running a suite, it clones a base CSV sink per experiment to generate separate CSV files per experiment ¹⁰⁹ ¹¹⁰. This sink likely also sanitizes any dangerous content (like Excel formulas) by default – in the code when cloning, it copies `sanitize_formulas` and `sanitize_guard` flags ¹¹¹. That indicates the CSV sink has logic to prepend apostrophes to formula-like entries to prevent Excel injection.
 - `ExcelSink`: generates an Excel workbook (probably using `openpyxl`), which is included in optional deps ¹¹²).
 - `VisualReportSink`: produces visual charts (using matplotlib/seaborn).
 - `LocalBundleSink`: maybe collects all output files into a directory or archive.
 - `ZipBundleSink`: zips multiple artifacts into one file for easier distribution.
 - `SignedSink`: computes a cryptographic signature for a bundle of artifacts. This likely interacts with the `core.security.signing` module. The signing process might generate a manifest of files with hashes and sign that using a private key (the details depend on implementation; it could use an HMAC or RSA signature).
 - `RepositorySink`: as discussed, pushes results to a remote repository via API. Each sink can also declare what artifacts it **produces** or **consumes** (via an `artifacts` config). For instance, a `VisualReportSink` might produce a file artifact (the chart image) that a `ZipBundleSink` will consume. The artifact pipeline uses these declarations to order sink execution ¹¹³ ¹¹⁴. Sinks also carry a `_elispeth_security_level` attribute (set during instantiation from config) ⁴⁶ ¹¹⁵, which is used to enforce that they only receive data allowed for their clearance ³. Finally, each sink is typically instantiated by the orchestrator at startup (global sinks from config) or by the suite runner per experiment (overriding or cloning sinks for each experiment as needed) ⁴¹ ¹¹⁶. The actual writing is coordinated at the end of run by the artifact pipeline, which calls `sink.write()` in the correct order ¹¹⁷ ¹¹⁸.

Given these layers, the relationships are largely **hierarchical** with some cross-cutting: - The **CLI** invokes the **Orchestrator/SuiteRunner**, passing in the config-derived **Settings** and flags ²⁵ ⁴⁰. - The **Orchestrator** uses **Plugin Registries** to create concrete plugin instances (data source, LLM client, etc.) ⁷³ ⁷⁴. These instances adhere to abstract base classes (interfaces defined in `core.protocols` or

`plugins.*.protocols`) and are stored in orchestrator or runner fields. - The **ExperimentRunner** (instantiated by orchestrator) holds references to the LLM client, lists of plugin instances, and sink list ^{74 119}. It then orchestrates their use: for each data row, it calls the LLM client, then plugins, accumulates results, and after looping, calls aggregator plugins and then triggers sink output via the artifact pipeline. - **Dependency injection and inversion:** Note that orchestrator injects dependencies into runner (instead of runner querying global state). The runner doesn't directly know about orchestrator or CLI; it just knows about the LLM client and plugins given to it. Similarly, plugins often don't know about each other (except via the context and the data passed in). This reduces coupling – e.g., a row plugin just processes a row and doesn't need to know where the data came from or where it's going. - **Component relationships in code** can be seen in import and instantiation patterns: - The CLI imports the orchestrator and suite runner classes ¹²⁰, indicating the flow of control goes from CLI to those core classes. - The orchestrator imports factory functions like `create_row_plugin`, `create_aggregation_plugin` from the plugin registry ¹²¹. This shows orchestrator's dependency on the plugin subsystem to instantiate plugins by name. - Orchestrator also imports `ExperimentRunner` and creates it ^{122 123}, meaning the ExperimentRunner is a subordinate component handling detailed logic at runtime. - The plugin registry likely registers concrete plugin classes in dictionaries keyed by name (the test uses `register_row_plugin("single_run_row_plugin", factory)` to add a custom plugin at runtime ¹⁴, which implies internally `plugin_registry` stores that). Then `create_row_plugin(defn, parent_context)` looks up by name and calls the associated factory to get an instance ⁷³. This design decouples the orchestrator from having to know plugin class definitions – it only deals with names and definitions. - Sinks are created either by `sink_registry.create(name, options)` as seen in `ExperimentSuiteRunner._instantiate_sinks()` ¹²⁴ or passed pre-made from orchestrator config. The artifact pipeline then links sinks together using their produce/consume declarations.

Cross-cutting concerns in the logical view: - **Logging:** The use of Python's logging is pervasive but simple – e.g., `logger.info("Running single experiment")` in CLI ²⁵ or logging warnings from validation results ³⁵. The system doesn't have a dedicated logging service or structured log aggregator; it relies on the execution environment to handle log capture. However, via middleware, one could integrate an enterprise logging (for instance, an `AuditLoggingMiddleware` could send structured logs to a SIEM). Indeed, presence of an `audit-logging.md` doc suggests such patterns are considered even if not hard-coded. - **Security context propagation:** Through the `PluginContext` and security level normalization, the orchestrator ensures every plugin and sink knows the classification of data it handles ^{48 125}. This is cross-cutting in that it spans the data source (which might have a `security_level` attribute) to the LLM client and outputs. The code calls `resolve_security_level()` combining e.g. datasource and llm's own levels ⁴⁸. Then it tags contexts and sinks. This means at any point, a plugin could query its attached context to know security level or determinism constraints. The artifact pipeline enforces those constraints globally when linking artifacts ³. - **Error handling:** Largely concentrated in the runner loop – catching errors per row and storing them, etc. The design choice is not to use exceptions to unwind the entire orchestrator on a single row's error, but rather to capture it and continue (unless an early-stop triggers a halt). This makes the system robust to individual data issues or API hiccups. - **Configurability:** The logical architecture heavily relies on configuration-driven behavior. The `OrchestratorConfig` dataclass holds dozens of fields that mirror YAML config sections ^{126 127}. Because of this, many components get their instructions from that config rather than hardcoded decisions. For instance, concurrency is off by default unless `concurrency_config.enabled=true` in YAML, at which point the runner will use threads if row count threshold is met ⁷. Similarly, if no `early_stop` is in config, none is applied; if present, the system will instantiate the plugin to enforce it. This pattern means the logical layer is quite dynamic – it assembles the exact pipeline for a run based on config. The code for merging config layers in suites ^{128 129} shows how it programmatically builds the effective configuration for each experiment. This

approach ensures consistency (each experiment runner is configured via the same process), but also means thorough config validation is needed to catch mistakes (the code uses `validate_settings()` before running, which checks the YAML against expected schema) ²⁶.

In **C4 model terms**, one could consider: - *Components* like Orchestrator, Runner, SuiteRunner as the core “brain” components. - *Containers* not in the Docker sense but as groupings: the CLI and these orchestrator components run in the same process container. - External systems (OpenAI, Azure, etc.) would be separate boxes the system interacts with via well-defined interfaces.

From a logical view, the **flow of control** for a typical single experiment run is: 1. **CLI** parses args, loads config into a `Settings` object (which contains a pre-instantiated `datasource`, `llm`, `sinks`, and an `OrchestratorConfig` among other things) ^{130 131}. 2. CLI calls `ExperimentOrchestrator(datasource, llm_client, sinks, config, rate_limiter, cost_tracker, ...)` to initialize orchestrator ²⁵. 3. Orchestrator's `__init__` creates plugin instances for all plugins defined in `OrchestratorConfig` (row, agg, validation, early-stop) using registries ^{73 132}, and sets up the `ExperimentRunner` with those plus other configs ^{74 119}. 4. CLI then calls `orchestrator.run()`. Inside run: - Load `DataFrame` from `datasource` ⁸⁵. - Apply any row limit (`max_rows`) ¹³³. - Pass `DataFrame` to `runner.run(df)` and receive a `payload` dict ¹³⁴. 5. The `ExperimentRunner.run(df)` executes: - Prepares prompt templates via `PromptEngine` (compiling Jinja templates with defaults) ^{135 136}. - Iterates over `DataFrame` rows, either sequentially or in parallel threads depending on concurrency settings ^{137 138}. For each row: * Renders system and user prompt by filling in row data (the `_render_prompts` step uses compiled templates and the context) ^{139 140}. * Calls LLM client's `generate()` with the prompts (through `_execute_llm()`), which handles retries and middleware) ^{141 142}. * Receives LLM response (or error). If error after retries, creates a failure record with error info ¹⁴³. * If successful, attaches retry metadata to the result and calls each row plugin's `process_row` to add metrics to the record ¹⁴⁴. * Marks the record with security level (so each record carries that classification) ⁸⁷. * If no failure, adds record to results list; if failure, adds to failures list. The `handle_success` callback also appends to a checkpoint file if configured. * Checks early-stop plugins after each success (maybe the plugin says “stop now”), and if triggered, breaks the loop ^{145 146}. - After looping, sort results back to original order, then call each aggregator plugin's `finalize(results)` to get aggregates ⁹². - Build a `payload` dict with `"results"`, `"failures"` (if any), `"aggregates"` (if any) ⁴⁴. - Compute a metadata summary: number of rows, retry statistics, attach aggregates and cost summary if present, and the overall security & determinism level for the dataset ^{147 148}. - If an early-stop occurred, include the reason in both payload and metadata ²². - Then create an `ArtifactPipeline` with all sinks bindings and call `pipeline.execute(payload, metadata)` ¹¹⁷. This will orchestrate calls to each sink's `write()`: essentially distributing the `payload` and parts of it to sinks based on what they consume/produce. For example, the CSV sink likely consumes all results (it writes them), whereas a signing sink might consume the list of artifacts produced by earlier sinks (the pipeline collects artifacts as sinks run). - Once sinks have executed, the payload (with possibly additional entries like signed hashes if the signing sink added them) is returned. 6. Back in CLI, for a single run, it converts the results into a pandas `DataFrame` and prints a preview or writes to a CSV if requested ^{149 70}. It also logs any row-level failures that happened as error logs for visibility ³⁶. 7. For a suite run (if multiple experiments), the `ExperimentSuiteRunner.run(df, defaults)` process differs in step 4: - SuiteRunner loops through each experiment config, instantiates sinks for it (either using global ones or overriding) ¹¹⁶, calls `build_runner` to create a configured `ExperimentRunner` for that experiment ^{78 79}, triggers its `.run(df)` similar to above ^{150 151}, and collects payloads per experiment. - It calls `on_experiment_start` and `on_experiment_complete` on any LLM middleware for each experiment ^{107 105}. - After all, it handles baseline comparisons as described, and then calls `on_suite_complete` on middleware ¹⁵². - Returns a dictionary of results per experiment.

The above logical flow highlights the layered responsibilities: - CLI: parse and present. - Orchestrator/SuiteRunner: configure and iterate experiments. - Runner: iterate rows and apply plugins and controls. - Plugins: encapsulate specialized processing at various points. - Pipeline: finalize artifact delivery to outputs.

The code structure strongly reflects this separation of concerns, which is a positive sign for maintainability. For instance, if one wanted to add a new type of plugin (say a “column transformation plugin” that preprocesses the DataFrame), they could do so by introducing it in config and updating the runner to apply it at the right point, without affecting CLI or orchestrator logic significantly.

In conclusion, Elspeth's logical architecture is a **plugin-oriented pipeline** governed by a central orchestrator. The layering ensures that the core logic (prompting LLMs and capturing results) remains mostly unchanged regardless of which data source or sink or metric is in use – those are all plug-and-play. This modular layering is intentional to support the “**Composable plugin system**” touted in the README ⁶. Evidence of that modularity is seen in how the orchestrator references plugin interfaces and not concrete classes, and how easily Dummy plugins were injected in tests to simulate the system behavior ¹⁵³ ⁸⁶.

2.2 Physical Architecture (Deployment View)

The physical view describes how and where the system is deployed, including processes, containers, and network topology. Elspeth, being essentially a Python program, has a relatively simple deployment footprint by default: - It runs as a **single process on a single host** when invoked. This can be a developer's workstation, a server in a data center, or a container in the cloud. There is no persistent server or multi-process cluster built into the system. Each invocation is ephemeral – the process starts, does the work, and terminates. - There is no built-in clustering or scaling at the process level (scaling would be achieved by running multiple processes in parallel if needed, orchestrated externally by the user or an automation tool). - **Memory and CPU usage** depend on the dataset size and concurrency settings. With concurrency enabled (`max_workers>1`), Elspeth will spawn threads within the process (via `ThreadPoolExecutor`) ¹⁵⁴. These threads share memory (the pandas DataFrame, etc.) and are managed by Python's runtime. The code uses threads rather than separate processes, likely because the tasks involve I/O waits (API calls) or because using threads allows shared state like rate limiter counters without explicit IPC. This means Elspeth's parallelism is constrained by Python's GIL if tasks are CPU-bound, but in this case, the waiting on network calls is the main time sink, so multithreading yields concurrency. - If extremely large datasets are used, memory could become a factor (since the entire DataFrame is loaded into memory by `datasource.load()` and possibly duplicated when converting to lists of dicts or DataFrames for output). The code does allow limiting rows (`max_rows` config) ¹⁵⁵ ¹³³, which can be used to avoid OOM in memory-limited environments. - **Processes and Services:** From deployment perspective, Elspeth might be packaged as a Python package installed in a virtual environment. The provided Makefile and bootstrap script indicate a typical installation (creating `.venv`, installing dependencies) ¹⁵⁶. There's no mention of a dedicated daemon or OS service – it's run on demand. - If integrating into automation, one might run Elspeth via a CI/CD pipeline or a scheduling system like Airflow by calling the CLI.

While Elspeth itself is one process, it **interacts with external infrastructure**: - **Networking:** All network interactions are client-initiated (outbound). The code will open HTTPS connections to: - OpenAI or Azure endpoints (usually `api.openai.com` or an Azure endpoint domain) for LLM queries. - Possibly Azure Blob endpoints (`*.blob.core.windows.net`) if using cloud storage. - GitHub or Azure DevOps (likely `api.github.com` or an Azure DevOps REST endpoint) if those sinks are used. - Azure Content Safety or other Azure services if corresponding middleware is enabled (for example, if an Azure content policy middleware calls an Azure endpoint). - These external calls might need proxies or custom

certificate authorities in enterprise networks. The docs hint to “*enable corporate CA bundles and proxy settings where required*” when using requests ⁴⁵, which is a deployment consideration: e.g., setting environment variables for HTTPS proxies or using `REQUESTS_CA_BUNDLE`. - There is no incoming network port that Elspeth listens on, so no firewall rules are needed to allow inbound traffic to Elspeth specifically. Standard egress firewall rules must allow the above outbound calls if those features are used.

- **Containers:** Although not explicitly provided, Elspeth could be containerized. The absence of a Dockerfile in the repo suggests it’s not officially containerized yet. But since it relies on Python and some native libraries (pandas, numpy, etc.), a container image would need those. In a deployment scenario for ATO, one might create a Docker image with Python 3.12, install Elspeth and its dependencies, and use that for consistent deployment. The physical architecture in that case is Elspeth container(s) running on a container orchestration platform (Kubernetes, ECS, etc.) on demand. The stateless nature of each run means you could spin up multiple containers to parallelize experiments (different experiments or different segments of data, as orchestrated externally).
- **Resource Isolation:** Because each Elspeth run can connect to external APIs, one might run it in a restricted environment if needed. For instance, an accredited environment might require running the process on a VM that has no internet access unless explicitly whitelisted. In that case, using the `mock` LLM plugin would be necessary for offline runs (or connecting to an internal LLM endpoint). The code supports a fully offline mode with no external calls (as proven by tests using DummyLLM and local CSV).
- **File System:** Elspeth reads config and possibly local data files from disk, and writes output files to disk. Therefore, it requires access to the file system (read/write in certain paths). By default, outputs go under `outputs/` directory (as per README Quick Start) ¹⁵⁷. In a container, this might map to a volume or artifact storage. If using Azure blob for outputs, it writes via API rather than local disk. But intermediate outputs (like the combined results DataFrame for preview, or the coverage XML for tests) are local.
- **Physical Security:** If the data is sensitive, the environment should ensure disk encryption and proper access controls around where the process runs. This is outside the code, but important in deployment.

Deployment Scenarios: 1. *Developer laptop:* Run `make sample-suite` as in README ¹⁵⁸. This downloads a sample CSV, runs experiments with a mock LLM, and writes reports locally. The architecture is simply Python process + local file system. 2. *On-prem server for nightly experiments:* The process might be scheduled (cron or via CI) to run with access to internal data sources. It might use the OpenAI API via a proxy. The server needs connectivity to those APIs, and secrets (API keys) would be provided via environment or a secure store and loaded into config at runtime. 3. *Cloud deployment in an enterprise:* Possibly running in a secure container with no internet – in which case, any references to external LLM must be replaced by either an internal model endpoint or the run must use recorded data. Or one could run it in an approved cloud environment with internet but restricted to certain API endpoints (with monitoring). 4. *Integration with other tools:* Because it’s CLI, one could incorporate Elspeth into larger workflows (like calling it from a Jupyter notebook or a web service). If someone wraps the CLI call inside a web service (for example, a Flask app that triggers an Elspeth run on request), that would create a new boundary (the web service to CLI call). But that’s an extension beyond what the codebase itself provides.

Network architecture within Elspeth is straightforward: no internal network of services, just calls from the main process to external host addresses. The concurrency is handled in-process; no service mesh or

inter-process comm. The environment likely has to handle secrets – e.g., Azure Managed Identity credentials come from the host (if running on an Azure VM or container with MSI, azure-identity library will fetch a token from the Azure Instance Metadata Service at `http://169.254.169.254/...`—this is one notable network call to mention. The azure-identity default behavior in a cloud environment is to contact the IMDS endpoint for a token, which is a link-local call from the VM. That is another kind of integration: connecting to a metadata service for auth. In a physical view, that means if running on Azure, the VM needs that IMDS access. If running off-cloud, azure-identity will fall back to environment client ID/secret if provided).

Infrastructure as Code / Config: There's no Terraform or K8s manifest in the repo. The expectation is that the user will handle deployment. For ATO, describing a typical deployment might be needed: - Possibly, an **ATO would consider a hardened VM or container** where Elspeth runs. E.g., a container built on a FIPS-enabled base image if required (especially since cryptography might be used for signing). - Also consider **scaling**: If the workload is heavy, one might run multiple Elspeth processes in parallel on different input partitions (since one process uses only one CPU effectively when not waiting on I/O, aside from multi-threading up to 4 threads by default). - The **physical data flow**: Input data might reside in a secure data lake or in Azure Blob; Elspeth pulls it to memory, processes, and might store outputs back to another secure location. The data might not persist on the local disk unless explicitly configured to (except ephemeral working files). - **Statefulness**: Each run is stateless in the sense that it doesn't rely on data from previous runs (aside from optionally reading a checkpoint file if resuming). If checkpoint is used, that file is a small JSONL file listing processed IDs ¹³, which should be stored in a location safe from tampering if it's used to skip records (to prevent replay or skipping unauthorized). But that's a minor detail – likely in practice, checkpoint file would be local or on a mounted storage.

- **Security aspects physically:** ensure that any credentials (like the OpenAI API key or Azure storage key if used directly) are provided via environment variables or secure config that is not world-readable. The code does not print secrets (no such logging observed), and any error would ideally not dump secrets (unless an exception deep in azure libs might, but presumably not). This is part of deployment hardening – use Azure Key Vault or environment injection for sensitive config fields (the `docs/compliance/configuration-security.md` likely covers recommending not to put plaintext secrets in YAML, or to use placeholders).

Overall, the physical architecture is **lightweight** – one can deploy Elspeth on existing compute without needing additional service infrastructure (databases, caches, etc.). This simplicity can be an advantage for ATO as there are fewer moving parts to secure (just the runtime environment and outbound connections).

For completeness, one could draw a simple deployment diagram: a box for Elspeth CLI app on a host, arrows going out to “LLM API”, “Azure Blob Storage”, “GitHub API” as needed. Internally, the host has Python environment, and maybe an arrow from host to “Key Vault (credentials)” if that's used, or to “Metadata Service” for identity.

Without a container or microservice architecture, aspects like *load balancers*, *service discovery*, etc., are not applicable. The “Physical” architecture is essentially the runtime environment plus integration points.

2.3 Data Architecture (Data Flow and Model View)

In Elspeth, data flows through the system in well-defined stages. The primary data elements include: - **Configuration data**: defines what to do (experiments, prompts, plugins, etc.). - **Experiment input**

data: the actual dataset on which experiments run (tabular records). - **Intermediate data:** prompt templates, generated prompt strings, LLM responses (content and metadata), metrics. - **Output data/artifacts:** final results, aggregated metrics, reports, and possibly derived files (charts, signed bundles).

Data Models / Structures: - The main in-memory data structure for input is a **pandas DataFrame**. When `datasource.load()` is called, it returns a DataFrame with columns corresponding to the data fields (for example, columns like `APPID`, `name` as seen in tests ¹⁵⁹). Each row represents one unit of experiment (one prompt fill). - DataFrame may have attached metadata: for instance, in code we see `df.attrs.get("schema")` ⁹⁶ and `df.attrs.get("security_level")` ¹⁶⁰. This suggests that DataSource plugins can attach a schema (perhaps a JSON Schema or custom object) describing the data, and a security level classification for the entire dataset. The schema could include column types or allowed ranges, used by validation plugins to ensure e.g. prompts are compatible. - The **schema validation** process, hinted by `validate_schema_compatibility` function call ¹⁶¹, likely ensures that if a plugin expects certain columns (like a particular metrics plugin might require a column "score"), the data provides it. If not, it could warn or error (depending on config: note `on_schema_violation` attribute set to "abort" by default ¹⁶² – meaning by default, if schema validation fails, it aborts the run).

- Once the DataFrame is loaded, it is not directly fed to LLM calls; first it's possibly filtered by `max_rows` (limiting how many rows will be processed) ¹³³, and then each row is converted into a **context dictionary** for prompt rendering. The function `prepare_prompt_context(row, include_fields=prompt_fields)` does this conversion ¹⁶³. Essentially, it takes the pandas Series for a row and produces a Python dict (probably `{column_name: value}` for each column, but possibly filtering to only the fields listed in `prompt_fields` config if provided). This context dict is the data used to fill in the prompt template.

- **Prompt Templates:** Defined by the combination of the system prompt string and user prompt template (with placeholders). These are compiled into `PromptTemplate` objects which hold:

- `template` (the compiled Jinja2 template object),
- `required_fields` (fields that must be present to render),
- `defaults` (default values for some fields if not in context),

- metadata (like template name). Each experiment has at least one user prompt template (the primary one), plus a system prompt. Additionally, if `criteria` are defined (alternative prompts for the same input), each criterion has its own template (which if not explicitly provided might fall back to the main user prompt text) ¹⁶⁴. The PromptEngine compiles each criterion template as well, storing them in a dict keyed by criterion name ¹³⁶. Internally, these templates ensure that if a field is missing, a `PromptRenderingError` is raised rather than silently failing ¹⁶⁵. This ties into data validity: if the context dict doesn't have a needed field, that row's processing stops with an error (which will be recorded as a failure) ⁹⁵.

- **LLM Request and Response:** The data exchanged with LLMs is structured:

- The code creates an `LLMRequest` object for each attempt (this likely holds `system_prompt`, `user_prompt`, and `metadata`) ¹⁴¹. The metadata usually includes `attempt` count and possibly experiment name or other context ¹⁶⁶. Middleware can augment this request object (for example, injecting a user ID or removing sensitive content).
- The LLM plugin returns a **response** as a Python dictionary. By convention, responses contain at least a `"content"` field with the text output. They may also include `"prompt"` (like the test dummy does) or other metadata from the API, and a `"metrics"` sub-dictionary for any metrics

(like token usage, model confidence if any, etc.). For example, the OpenAI library typically returns an object with choices; perhaps the plugin normalizes it to `{"content": "...", "usage": {...}}`. We see in the code that after receiving a response, the cost tracker might add metrics to it (like tokens and cost) under `response["metrics"]` ¹⁶⁷. Also, the retry logic ensures that the response dict contains a `"retry"` field with attempts info ¹⁶⁸. Indeed, the runner later uses that to compute `retry_summary` ²¹.

- If multiple LLM calls are made for one input (because of multiple criteria prompts), the code collects them: `_collect_responses` might produce a dict of responses for each criterion plus a primary response ¹⁶⁹. The test `test_orchestrator_with_criteria` shows that in the payload, `result["responses"]` contains keys for each criterion name mapping to its response ¹⁷⁰. And `result["response"]` for the primary response is also present ¹⁷¹. So the data model for a single record's outputs can be:

- `"response"`: the main response dict,
- `"responses"`: a dict of `criterion_name -> response dict` (for alternate prompts),
- `"metrics"`: a dict of any metrics from row plugins,
- `"row"`: the input context (sometimes stored for traceability, as seen in failure case they store `failure["row"] = context` ¹⁷²).
- `"retry"`: if that row's request had retries, an object with attempts and history.

- **Result Records and Aggregation:** After processing all rows, the runner compiles a list of **result records**, each as described above, and a list of **failures** (each a dict with at least `"row": context` and `"error": message`). It then applies aggregator plugins:

- The aggregator plugins return a derived value (could be scalar or dict). The runner stores these in an `"aggregates"` dict keyed by plugin name ⁹². For example, aggregator `"single_run_agg_plugin"` produced `{"count": 1}` in the test, which appears in `payload["aggregates"]` ⁸⁹.
- Baseline comparisons produce another layer of aggregated data: for each experiment (except baseline), a `"baseline_comparison"` dict may be added, mapping baseline plugin names to their output (which could be simple metrics or structured diff). We see code collecting `comparisons[plugin.name] = diff` and adding that to payload ¹⁰¹.

• Data Persistence and Formats:

- In-memory, data is mostly Python dicts and pandas DataFrames.
- When writing to disk or external, the format depends on the sink:
 - CSV sink writes a CSV file (UTF-8 text). Each row becomes a line, each metric possibly becomes a column (the CLI flattens metrics into columns using `_flatten_value` for nested metrics before writing CSV) ¹⁷³ ¹⁷⁴. For example, if a response has `{"metrics": {"score": 0.5}}`, CLI's `_result_to_row` will flatten it to a column `metric_score` in the CSV ¹⁷⁵.
 - Excel sink uses `openpyxl` to produce `.xlsx` (binary Excel format).
 - JSON output (if any sink or export uses JSON) – possibly the CLI allows exporting the suite configuration to JSON/YAML via `--export-suite-config` ¹⁷⁶, which would output the merged config. But result-wise, not sure if there's a JSON result sink; perhaps not explicitly except the baseline Signed sink might produce a JSON manifest or similar.
 - Markdown/HTML: Possibly a sink to generate Markdown or HTML summary (the `VisualReport` might output an HTML report with embedded images).

- In-memory results (for programmatic use): The CLI converts results to pandas DataFrame for quick display ¹⁷⁷, meaning the result records (list of dicts) are readily tabular. Each dict's keys become DataFrame columns (which means nested structures like metrics and responses are flattened by `_result_to_row` as mentioned).
- **Artifact Pipeline:** This component ensures that if one sink produces a file artifact that another sink needs to include, it will run them in correct order and pass the artifact references. The data model here includes:

- `ArtifactDescriptor` (with fields: name, type, optional alias, security_level, persist flag) ¹⁷⁸ – defines an artifact a sink produces or needs.
- `Artifact` actual objects stored in `ArtifactStore` with an id, security level, possibly content or reference to file.
- `ArtifactRequest` (token+mode) to express a sink's consumption needs (e.g., token "csv" with mode "single" means "give me the CSV artifact") ¹⁷⁹ ¹⁸⁰. The pipeline uses these to build a dependency graph of sinks: For example, a `SignedSink` might declare it **consumes** all artifacts (`{"token": "all", "mode": "all"}`) and **produces** a signature file. The pipeline will ensure all other sinks produce their files first, then the `SignedSink` runs. If a sink's declared `security_level` is lower than an artifact's level, pipeline will raise an error (as seen in code) ¹⁸¹, aborting execution to prevent a data leak. This means the data architecture has an implicit **type system for artifacts** (maybe types like "csv", "excel", "image", etc.). `validate_artifact_type` is called to ensure tokens are recognized types ¹⁸².

• Indices and Keys:

- There isn't a relational database, but the DataFrame typically has an index. In code they iterate with `df.iterrows()` which gives an index and row ¹⁶³. They use their own sequential `idx` counter as well. If the data has a unique identifier column (like "APPID" in tests), it can be used as a stable key. In fact, the checkpoint logic uses a checkpoint field (default "APPID") to mark processed rows ¹³. This implies the input data should have a primary key column for checkpointing to work reliably (if not, one might use index but index can shift).
- The context `row_id = context.get(checkpoint_field)` ¹⁸³ is fetched to identify the row in checkpoint logs. If `processed_ids` contains it, the row is skipped (so we don't repeat already done work) ¹⁸⁴.
- This is the only use of a "key" concept we see – essentially, "APPID" or any configured field acts as a unique row identifier for resuming runs.

• Data transformations:

- The transformation from raw input row to prompt (via template) is a big one. The Jinja environment is configured not to autoescape (since typical content is not HTML) ¹⁸⁵. It uses `StrictUndefined`, so missing fields cause an exception rather than default to blank ⁴. There's also `_auto_convert` that automatically changes `{field}` style placeholders in prompt strings to Jinja's `{{ field }}` format for compatibility ¹⁸⁶. This lets users write simpler templates without braces doubling, while still using Jinja under the hood.
- If prompt aliases or prompt field mappings are configured (there is `prompt_aliases` in `OrchestratorConfig` ¹²⁶), those might transform the context keys (e.g., aliasing a column name to a shorter token in template). The code doesn't show where `prompt_aliases` is used, but likely `PromptEngine` or `prepare_prompt_context` might apply it (not in the snippet though).

- **Metrics merging:** The runner merges various metrics from different sources. For example, if the LLM response has its own `"metrics"` (like token counts), and row plugins return metrics, and cost tracker adds cost metrics, they all end up in the final record's `"metrics"`. The code ensures to merge these properly. There is `_merge_response_metrics` method referenced in runner (not shown in snippet) that likely merges metrics from primary and criteria responses into the main record ¹⁸⁷.
- At output time, the CLI's `_result_to_row` flattens nested metrics and responses for tabular output ¹⁸⁸ ¹⁸⁹ (e.g., it will create columns like `llm_content` for main response content, and `llm_{name}` for each named response content, and separate columns for each metric from each plugin).
- **Caching:** There is no long-term caching of results between runs in code (besides checkpointing). However, within a run, certain things are cached:
 - `ExperimentSuiteRunner` caches shared LLM middleware instances so that if multiple experiments use identical middleware config, it reuses one instance instead of creating many ⁸². This prevents duplication of e.g. a heavy telemetry client across experiments. The key for caching includes the config and security level to ensure same config yields same instance ¹⁹⁰.
 - `PromptEngine` might cache compiled templates by holding them in the `PromptTemplate` objects – repeated rendering of the same template is cheap after compile.
 - No explicit memoization or reuse of LLM responses – if the same input appears twice, it will call the LLM twice (unless user manually deduplicated input or used checkpoint to skip).
 - If needed, one could implement a caching plugin (e.g., to cache LLM responses by prompt to avoid re-querying identical prompts), but that's not present by default.
- **Data Volume:**
 - Input size: Because the whole dataset is loaded into a `DataFrame`, data volume is limited by memory. If an input CSV is huge (say millions of rows), this could be an issue. The design does not chunk streaming data – though interestingly, a roadmap file named "streaming-datasource-architecture.md" is in the repository, suggesting plans to handle larger-than-memory data by streaming or incremental processing. In the current implementation, each run is all-or-nothing in memory.
 - Output size: If generating an output for each input row plus metrics, the results can be as large as input in row count, with additional columns. Writing to CSV or Excel can handle fairly large outputs (Excel has limits ~1M rows). For extremely large outputs, a CSV might be tens of MBs or more. If needed, using a database sink plugin could be an extension to handle volume, but not built-in.
 - Reports like visual charts usually aggregate to manageable size (a chart image).
 - The signing output likely adds minimal overhead (a signature file and possibly a manifest listing file hashes).
- **Consistency and Transactionality:**
 - There is no transactional database to worry about atomic commits. However, consider the scenario of writing multiple outputs: the artifact pipeline tries to produce a consistent set of artifacts. If one sink fails midway (e.g., network error uploading to GitHub), what happens? The pipeline likely propagates the exception (since we see it raises errors for security mismatches;

presumably if a sink's write raises, it stops). That would mean some outputs may have been written (e.g., local CSV file might be already written) while others not. In an ATO context, one might want either all or nothing for outputs – currently, the design doesn't roll back already written files if a later sink fails. It simply reports failure in the run metadata (failures list or an exception if not caught). The user could rerun or manually handle partial outputs.

- Because each sink is independent (except for declared dependencies), partial results are possible. For example, CSV and Excel sinks might succeed, but a repository sink fails – you'd have local files but nothing in repo. The run's logs would show an error for the repo sink. Mitigation might be to run repository sink first (but they deliberately run it last with dry-run by default to avoid pushing partial info).
- Checkpointing is the only persistence between runs: it is appended to after each row success. The checkpoint file (JSON Lines containing IDs) is thus always consistent in that if a run aborts, all lines in it correspond to completed rows. On resume, those rows are skipped. If the run aborted in the middle of a row (unlikely since processing of a row is not interruptible by design, except by early stop which is triggered right after finishing a row), then that row wouldn't be marked processed and will be retried. So checkpoint ensures eventually all rows get processed across runs, albeit some might be reprocessed once if failure happened just after they completed but before writing checkpoint (which should not happen because they append to checkpoint as part of `handle_success`, before any potential heavy operations afterwards).

To illustrate the **data flow end-to-end**: 1. **Configuration** (YAML) is read and validated. This yields: - `Settings.datasource` (object to load data), - `Settings.llm` (LLM client object), - Lists of plugin definitions (row, agg, etc. still as dicts), - `Settings.sinks` (list of result sink objects pre-instantiated), - `OrchestratorConfig` (with prompt templates, plugin definitions, etc. in structured form). 2. **Data Ingestion**: `datasource.load()` is called → returns `DataFrame df` ⁸⁵. The `DataFrame` might have `df.attrs['schema']` if data source set it (for example, an advanced data source might attach a schema read from a metadata file or infer types). 3. **Experiment Processing**: - Orchestrator (or SuiteRunner for each experiment) sets up an `ExperimentRunner` with references to: * The `DataFrame df`, * `prompt_system` and `prompt_template` strings (from config) ¹⁹¹, * `prompt_fields` list (to filter context, if given) ¹⁹², * `criteria` list (if multiple prompts per row) ¹⁹³, * `row_plugins` list, `aggregator_plugins` list, etc. * `llm_middlewarees` list, `rate_limiter`, `cost_tracker` objects. - The runner compiles the prompt template into `PromptTemplate` objects (and similarly for criteria) ¹⁹⁴ ¹³⁶. - Iteration: for each row in `df`: * Create `context` dict from row ¹⁶³. If `prompt_fields` is specified, only include those keys; if not, include all columns as keys. * If any `prompt_aliases` were configured, possibly rename context keys accordingly (not shown explicitly, but presumably before rendering). * Render `system_prompt` and `user_prompt` by plugging context into templates ¹⁹⁵. If a required field is missing, raise `PromptRenderingError` caught and added to failures ⁹⁵. * If criteria present: for each criterion, similarly render a criterion-specific user prompt variant. Keep them in a dict of `{criterion_name: rendered_prompt}`. * Formulate request(s) to LLM: - For primary prompt: fill `LLMRequest` with system + user prompt text and metadata (like row index, experiment name) ¹⁴¹. - For each criterion, possibly call LLM separately (unless the plugin supports multi-prompt in one call, but likely not; it's probably sequential calls). - Middleware's `before_request` may modify the request (e.g. injecting headers or filtering text) ¹⁰⁶. - Rate limiter's `acquire()` might block until allowed and then proceed within context manager (ensuring count is tracked) ⁹. - Call `llm_client.generate(system_prompt, user_prompt, metadata)` – actual API call or local generation ¹⁹⁶. - On receiving `response` dict: + Apply `after_response` middlewarees in reverse order ⁵⁷ (e.g., could remove some sensitive info or add tags). + Pass `response` and request to validation plugins `_run_validations(...)` ⁹⁴. If a validation fails, it might raise an exception, which would be caught and treated as a failure for this row (with error message, likely through `PromptValidationError` catch that returns a failure dict) ⁹⁵. + If `cost_tracker` is present, call `cost_tracker.record(response, metadata)` to extract cost metrics (like tokens used and

approximate cost in \$) ¹⁶⁷. This returns a metrics dict (e.g., {"tokens": 100, "cost_usd": 0.002}) which is then merged into `response["metrics"]` ¹⁶⁷. - If an exception occurred in `llm_client.generate` (network error, API exception), the code catches *any* Exception: + It logs attempt in `attempt_history` with status "error", error type, etc. If attempts remain, it waits then retries ¹⁹⁷ ¹⁹⁸. + If max attempts reached, it breaks out with `last_error`. The code then attaches the attempt history and attempt count to `last_error` as attributes ¹⁹⁹, and calls `_notify_retry_exhausted` (which might let middleware know all retries failed) ²⁰⁰. + It doesn't explicitly show raising the exception up, but since there's an `assert last_error is not None`, probably it will raise after that or return a failure object. Given how `_execute_llm` is used inside `_process_single_row`, likely if it fails, they return a failure dict instead of record. * After getting one or multiple LLM responses for the row: - Build a `record` dict with `"row": context` (for trace), `"response": primary_response`, and if multiple, `"responses": {name: resp_dict, ...}` ²⁰¹. If the LLM client itself returned metrics (like token counts or attempts), those are in `primary_response["metrics"]`. - Populate prompt metadata in the record (somewhere in `_populate_prompt_metadata`, which might store the actual rendered prompt text or parts of it if needed for debugging or reproducibility) ¹⁴⁰. - Attach the `"retry"` field from primary response (which has attempts info) into the record for easier access ¹⁴⁰. - Run **row plugins**: call each plugin's `process_row(record["row"], record["responses"] or record["response"])`. If a plugin returns a dict of metrics, merge those into `record["metrics"]`. The code does `self._apply_row_plugins(record, row_plugins)` which likely iterates and updates record ¹⁴⁴. The example plugin in test added `{"custom_metric": 7}` which ended up as `record["metrics"]` `["custom_metric"]=7` ²⁰² ⁸⁹. - Apply **security** **label**: call `self._apply_security_level(record)` which might tag the record with `record["security_level"]=active_security_level` (the code suggests they do this before returning record) ⁸⁷. Indeed, we saw in CLI `_result_to_row` they check `if "security_level" in record: row["security_level"]=record["security_level"]` ²⁰³, so each output row carries its classification. - At this point, if all goes well, we have a `record` with combined info. If any exception occurred in row plugins or during prompt rendering that's not caught, the broad `except Exception` in `_process_single_row` will catch it and produce a failure dict containing the error string and the context ²⁰⁴. * The runner then either returns a `record` (and no failure) or returns `None` and a `failure` for this row ⁹⁵. * The `handle_success` callback appends the (index, record) to results list, and writes checkpoint if enabled. The `handle_failure` appends the failure dict to failures list. * Early-stop check after each success: the runner passes the just-processed `record` to `early_stop` plugins: + If any returns a reason dict, it sets `_early_stop_reason` and signals the event ²⁰⁵ ²⁰⁶. The reason might include which plugin triggered and any details. + When event is set, the thread loop or for loop will break out soon after (they check at top of next iteration or in thread worker loop). - After processing rows (or breaking early), sort results back to original input order (since parallel threads may have collected out of order) ²⁰⁷. - Attach results list and failures list to payload. The payload structure emerging:

```
{
  "results": [ {record1}, {record2}, ... ],
  "failures": [ {failure1}, ... ], (only if any failures)
  "aggregates": { "plugin_name": value, ... }, (if any)
  "early_stop": { ... reason ... }, (if triggered)
  "cost_summary": { ... } (if cost_tracker was used and has summary)
  "metadata": { ... }
}
```

The `metadata` includes: * `"row_count"` equal to `len(results)` ²⁰⁸, * `"retry_summary"` if any retries happened ²¹, * `"aggregates"` copy, * `"cost_summary"` copy, * `"failures"` copy (so metadata carries failures as well, possibly for logging), * `"security_level"` and `"determinism_level"` which are the overall labels resolved for this experiment/data ¹⁴⁸, * `"early_stop"` reason if set ²². - The `ArtifactPipeline.execute(payload, metadata)` now runs: + It takes the configured sinks (each with `artifact_config` from when they were created). For each sink, a `SinkBinding` is prepared with: id (unique name + index) ²⁰⁹, plugin (sink type name) ²⁰⁹, the sink object, `artifact_config` (dict of produces/consumes from config), and its `security_level` ¹¹⁵. The `_prepare_binding` populates lists of `binding.produces` and `binding.consumes` by reading `artifact_config` and sink's own methods ²¹⁰ ²¹¹. For example, a CSV sink might produce `ArtifactDescriptor(name="results_csv", type="csv", ...)`, an Excel sink might produce type "excel", etc. A repository sink might not produce anything but consume e.g. "csv" (if it's meant to commit a CSV). + The pipeline then performs a **topological sort** of sinks based on dependencies: * If sink B consumes an artifact produced by sink A, ensure A runs before B ²¹². If there's a security level conflict (A's output classified higher than B's clearance), abort with error ³. * The pipeline likely uses a DAG approach: it repeatedly takes sinks whose dependencies are all satisfied, enqueues them for execution. + Execution: For each sink in order, it calls `sink.write(results, metadata)` (likely providing full payload or maybe filtered content based on what it consumes). Actually, it might not give the entire payload blindly: * The `ArtifactStore` might filter the payload for the sink. For example, a sink with `consumes: [{"token": "csv", "mode": "single"}]` expects one artifact of type "csv". The pipeline will find the artifact of type "csv" from the store and pass it (maybe as a file path or data) to the sink's `write` method or a parameter. * But since the sink objects are likely written to operate standalone (e.g., `CsvResultSink` knows how to get data from `results` list or `DataFrame`), possibly the pipeline's job is simpler: it doesn't filter payload for sink, but rather each sink knows what to do with the payload dictionary (like find results or metrics). However, the presence of produce/consume indicates a more advanced decoupling, but it may be future-proofing. In our analysis, we can assume each sink's `write()` is implemented to pick from the payload whatever it needs (commonly the "results" or "aggregates"). * After `sink.write`, if it produced any artifacts (e.g., created a file), the pipeline registers those in the `ArtifactStore` with a unique ID and type so downstream sinks can retrieve them ²¹³ ²¹⁴. + After all sinks, pipeline finishes. If any sink failed, an exception likely propagated (unless pipeline catches it and adds to failures; but code suggests it would raise, halting execution). - With sinks done, `ExperimentRunner.run()` returns the payload dict ¹¹⁷ (which now might have additional entries if sinks added something, though usually sinks don't modify payload except maybe adding file references or signature). 4. **Post-processing**: - The orchestrator's `run()` receives payload and simply returns it up ¹³⁴. - CLI receives payload. In single-run mode, CLI not only logs failures (iterating `payload.get("failures", [])` to log each error) ³⁶, but also converts the successful results to a `DataFrame` for preview or output. The code: * Builds list of `rows = [_result_to_row(result) for result in payload["results"]]` ²¹⁵. Each `result` is a dict with potentially nested metrics/responses, and `_result_to_row` flattens it: - It inlines top-level response content: e.g., `record["response"]["content"]` becomes `row["llm_content"]` ¹⁷³. - For multi-LLM responses, it inlines each named response's content under `llm_{name}` keys ²¹⁶. - It flattens metrics: any metric values (including nested ones) become columns like `metric_{metric_name}` ¹⁷³. For nested metric dicts, `_flatten_value` creates composite keys (joining nested keys with underscore) ²¹⁷. - It also carries over retry info (attempts counts) and the `security_level` if present ²¹⁸. * Then `pd.DataFrame(rows)` is created, which yields a tabular form of results ²¹⁹. * If `--output-csv` was given, it writes that `DataFrame` to a CSV file ³². If `--head` was requested (default 5), it prints the first few rows to console for the user ²²⁰. - In suite mode, the CLI doesn't itself convert to `DataFrame`; instead, it hands off to `SuiteReportGenerator` if `--reports-dir` was given ³³. `SuiteReportGenerator(suite, results).generate_all_reports(reports_dir)` is called. This likely uses the payloads of all

experiments plus the suite definition to create summary reports (like maybe a combined CSV of all experiments' results, charts comparing experiments, etc.). The code of `SuiteReportGenerator` (in `tools/reporting.py`) would define how cross-experiment data is aggregated for reporting. Possibly it uses pandas to combine each experiment's metrics, or it produces charts showing differences.

Entity-Relationship considerations: - The design is mostly procedural and does not define explicit classes for, say, "Experiment Result" separate from dictionaries. But conceptually: - One could consider an **Experiment** entity which has many **Result (RowResult)** entities and one **AggregateResult** (with metrics). - A **Suite** entity containing multiple Experiments and Baseline relationships. - However, these are not formalized as classes in code; they are implicit in how data is organized in payload structures.

Data Integrity & Validation: - Configuration is validated thoroughly by `validate_settings()`: it likely checks for missing sections, correct types, and possibly uses `jsonschema`. The code raises if any config errors ²²¹, and logs warnings for any deprecated fields etc. ²²². - The **input data** can optionally be validated: `--validate-schemas` CLI option triggers loading the DataFrame and then running schema compatibility checks without proceeding to LLM calls ²²³. In `_validate_schemas`, it likely uses the combination of data's schema and plugin expectations to ensure everything aligns ²²⁴. This mode ends early (return without running experiments) if validation is successful. It's a way to catch issues in data format ahead of time. - Each LLM response is validated by plugins (for example, if expecting JSON output, a validation plugin could attempt `json.loads` and error if fails). Or if using an allowed content list, a validation plugin could scan the text. - Results correctness is not automatically verified beyond what plugins check. Manual review of outputs is expected as part of experiment analysis.

ETL / Data pipelines: - Elspeth acts somewhat like an ETL pipeline for prompt experiments: Extract (from data source), Transform (via prompt + LLM → result), Load (into sinks). But each "transform" here is an LLM call, which is not a pure function of input data (it can be nondeterministic). The presence of `determinism_level` indicates they classify outcomes as "guaranteed" vs "unpredictable" – perhaps used to mark whether an experiment had any randomness (like if using temperature in LLM). - They allow deterministic offline replays by using a fixed random seed or the mock LLM (which always returns same outputs), which is critical for audit (ensuring reproducible runs when needed). This is part of data architecture insofar as they treat nondeterministic data carefully (e.g., capturing random seeds if needed – although not seen explicitly, they might consider it).

Message Queues / Streams: - The system does not use external message queues. If streaming results out was desired (like to stream partial results to a DB as they come), the architecture currently doesn't do that. It collects all in memory, then outputs at end. However, given Python's GIL, they probably wanted to avoid streaming in parallel to output as it might complicate concurrency. Instead, they ensure after all rows, everything is output in one go.

Summary: The data architecture is centered around the **pandas DataFrame** for input and lists of dictionaries for output, with JSON-like structures (dicts with nested fields for metrics and such) to carry rich information through the pipeline. The design chooses human-readable formats (CSV, JSON, YAML) at the boundaries to facilitate audit and integration. There is a conscious mapping of data fields to compliance constructs (like adding `security_level` field to records). The code's heavy use of dataclasses for config (ensuring type correctness) and structured logging of attempts and metrics all contribute to a system where data is both the input payload and the main output.

This architecture ensures that even though an LLM operation is not fully deterministic or structured, Elspeth wraps it in a lot of structure so that each piece of data (prompt input, output content, metrics

about the output, classification of the output) is captured and available for verification. This is key for an ATO context, where auditors want to see evidence of what data was sent and what came back.

2.4 Development View (Code Architecture & Build)

(Note: The template includes "2.4 Development View" and "2.5 Process View", which might overlap with what we've done, but I'll treat 2.4 as describing how the code is organized for developers and build processes, and 2.5 could be runtime process behavior which we've partly covered in 2.2 Physical but let's ensure we cover any missing development concerns.)

Code Organization: The project's source code is organized under `src/elspeth/` with subpackages reflecting functional domains: - `core/` - core mechanisms (orchestrator, runner, registries, security, config). - `plugins/` - all plugin implementations, further categorized (datasources, llms, experiments, outputs). - `datasources/` - possibly older or core data sources (e.g., `blob_store.py` might be a core utility class for Azure storage). - `tools/` - auxiliary tools like reporting generation. - `tests/` - comprehensive test suite covering each module.

This modular structure aligns with the logical architecture. For instance, `core/experiments/plugin_registry.py` holds the logic for registering and creating plugins by name ²²⁵. The presence of `core/controls` for cost tracking and rate limiting indicates those are considered part of core controls and not plugins (they are configured similarly but implemented in core). The code suggests that development considered backwards compatibility - e.g. docs mention some modules moved to `plugins/nodes/transforms/llm` etc. ²²⁶, implying a refactor where plugin locations changed but perhaps aliasing or notes ensure old references still work.

Build and Testing: - The project uses a `pyproject.toml` or similar for build configuration (since `requires.txt` is present, likely they use Poetry or pip for dependency management). - Running `make bootstrap` sets up environment and runs tests ¹⁵⁶. So tests are integral to the development process - all security and functional aspects have tests (including tests named `test_security_signing.py`, `test_validation_core.py`, etc. as seen in SOURCES list ²²⁷). - They use linters (ruff) and type checker (pytype) as part of linting via `make lint` ²²⁸, showing emphasis on code quality (static analysis). - The **Development view** likely covers how to extend the system: e.g., if adding a new plugin, one would implement the required interface in a new module under `plugins.*`, then register it in either `plugin_registry` or via an entry point. The plugin scaffolding might be eased by a script (we see a `scripts/plugin_scaffold.py` in search results, likely a helper to create a new plugin stub).

Quality attributes in development: - The code frequently uses Python's type hints and dataclasses (e.g., `OrchestratorConfig` is a dataclass ²²⁹, `Settings` is dataclass ²³⁰). This indicates a developer-focused clarity and some level of self-documentation (dataclass fields showing what config fields exist). - There are comments and docstrings explaining key classes (the orchestrator and runner have docstrings summarizing their role ²³¹ ²³²). This helps new developers or auditors understand intentions. - The existence of an `AGENTS.md`, `CLAUDE.md` etc. suggests they might have development notes possibly with AI agents for docs generation (since `.claude/agents` suggests they might be using an AI writing assistant ironically). - The dev docs in `docs/development/` cover upgrade strategy and dependency analysis, meaning the developers maintain documentation for future maintainers about how to safely update components and dependencies ²³³ ²³⁴.

Process View (runtime threads, concurrency): - We already detailed how concurrency uses threads. It might be worth summarizing separate from physical: - Within one run, a main thread orchestrates

overall flow. If concurrency enabled and row count is large, a thread pool with up to `max_workers` threads is created ¹⁵⁴. Each thread processes one row at a time using `worker()` function which calls `_process_single_row` ^{235 236}. - A lock is used to synchronize appending results to list so that there's no race in updating shared data structures ^{237 238}. - The rate limiter's `utilization()` is called by the main thread to decide if it should pause submitting new tasks to avoid overload ²³⁹. So the main thread might throttle dispatch to the pool if the external API is saturated. - Early stop uses an event; threads check `self._early_stop_event.is_set()` before doing heavy work and after finishing a row ^{240 241} to break out early. The main thread also checks the event between scheduling tasks ²⁴². - Because threads are used, any library calls that release the GIL (like I/O or network) truly run in parallel. The OpenAI API call (via requests) will release GIL while waiting for network, so multiple requests can be in flight concurrently. This is crucial for performance if calling many LLM requests. It's also why threads were chosen (ease of I/O concurrency). - No multi-process parallelism is used (like multiprocessing or distributed). This avoids complexities of data sharing and cross-process comm, which seems a deliberate design to keep things simpler for now. - The process lifetime ends after finishing output generation and returning from CLI main. There's no persistent state (except any files written).

Scalability note in development: - The code has parameters like `backlog_threshold` (if backlog of rows > 50, then parallelize) ²⁴³. This suggests they tuned it for a moderate number of parallel tasks vs overhead trade-off. If very few rows, they keep it sequential to avoid thread overhead. - There's also `utilization_pause` threshold in concurrency config (default 0.8, meaning if rate limiter utilization above 80%, pause new submissions) ²⁴⁴. This config helps prevent avalanche of tasks queueing if API is slower than we produce tasks. It's a nice adaptive strategy to mimic production load – an early sign of thoughtful architecture for realistic conditions.

Packaging: - The presence of `elspeth.egg-info` indicates it's installed as a package. For ATO submission, one might produce a wheel and a locked requirements set. The `dependency-analysis.md` serves as a guide to what dependencies are included, which is helpful for security review. It enumerates core and optional dependencies along with notes on CVEs or usage ^{68 245}. -

Licensing: The README says the License will be published and to contact maintainers in the meantime ²⁴⁶. For development, this means currently it's not open-sourced or the license is pending. That is a compliance risk because using it in production without a license could be problematic. They probably intend to finalize a license (maybe an enterprise wants to contribute back or something). - The code audit readiness: They maintain a `CONTROL_INVENTORY.md` and `TRACEABILITY_MATRIX.md` in docs, likely mapping code modules or config to NIST controls or organizational controls, which is part of the development artifacts for ATO. That means the developers integrated compliance requirements into the dev process, which is a strong sign of alignment between dev and security teams.

To sum up, the development view shows a codebase structured by feature, heavily tested, using modern Python features (dataclasses, type hints) and containing internal documentation to ease maintenance and audits. The build/test environment is straightforward with standard tools (pytest for testing, Makefile to automate tasks). There's an active approach to dependency management (explicit extras for optional features, monitoring vulnerabilities, and not bloating runtime with unused libs).

3. Component Catalog

This section provides detailed specifications of major components identified in the architecture. Each component is described by its purpose, responsibilities, key methods or interface, dependencies, and configuration. References to source code are provided to substantiate each point.

- **CLI Interface** (`elspeth.cli`) – *Purpose*: Parse user input and initiate experiment runs. *Responsibilities*: Command-line argument handling (using `argparse`), loading configuration, invoking validation or run modes, and basic result output formatting ^{24 32}. *Key Interface*: Executed via `python -m elspeth.cli` or console script. It provides usage help and options like `--settings <file>`, `--profile <name>` (to select config profile), `--single-run` vs default suite mode, toggles for metrics and live outputs, etc. ^{247 69}. After parsing, it calls `validate_settings()` and `load_settings()` to get a `Settings` object ²⁶, then either runs schema validation (`_validate_schemas`) or proceeds to execution (`_run_single` or `_run_suite`) ^{223 248}. It configures the root logger's level according to `--log-level` ²⁴⁹. *Dependencies*: The CLI imports the orchestrator, suite runner, and plugin classes it may need to instantiate (e.g., `CsvResultSink` for default CSV output) ¹²⁰. It depends on Pandas for preview formatting of results ²⁵⁰. *Configuration*: Accepts paths for settings and suite definitions, plus runtime flags. CLI does not maintain state between runs. *Notable Logic*: It handles creation of experiment templates and exporting suite config if requested (calls `create_experiment_template` and `export_suite_configuration` in `_handle_suite_management`) ^{251 252}. It also ensures that if `--live-outputs` is not set, repository sinks operate in dry-run (setting `dry_run=True` on them) ^{253 34}. *Error Handling*: The CLI catches configuration errors (via `settings_report.raise_if_errors()`) and will exit gracefully with an error message if config invalid ²²². Any runtime exceptions not caught inside orchestrator will bubble up (the CLI doesn't wrap the run in try/except, meaning a fatal error would crash the CLI and yield a traceback – in practice, orchestrator catches most operational issues).
- **Configuration Loader** (`elspeth.config`) – *Purpose*: Read and validate YAML configuration into structured objects. *Responsibilities*: Parsing the YAML file, merging profile sections, instantiating all required objects (datasource, LLM, sinks, controls) via registries, and producing a `Settings` dataclass that holds everything for execution ^{254 130}. It defines schema via helper dataclasses like `PromptConfiguration` and `PluginDefinitions` to organize config sections ^{255 256}. *Key Functions*: `load_settings(path, profile)` is the main entry – it loads YAML, selects the profile, and processes each part: datasource and llm via `_instantiate_plugin` and registries ^{130 67}, collects prompt and plugin definitions with `_collect_prompt_configuration` and `_collect_plugin_definitions` ^{257 258}, applies any prompt pack overrides (merging base config with a referenced pack of defaults) ^{259 260}, resolves sink definitions including injecting defaults for `security_level` and `determinism_level` if not provided ^{46 261}, and creates the `orchestrator_config` (`OrchestratorConfig` dataclass) populated with all these sections ^{262 263}. It also creates the `rate_limiter` and `cost_tracker` instances from their config if present ⁶⁷. *Dependencies*: Uses `PyYAML` for parsing YAML safely (`yaml.safe_load`) ²⁶⁴. Relies on registry singleton instances (`datasource_registry`, `llm_registry`, `sink_registry`) to map plugin names to classes ^{265 266}. *Configuration*: This component essentially translates user-provided config to internal config. It enforces defaults (if profile missing something, uses empty or default). It also calls `coalesce_security_level` and `coalesce_determinism_level` to validate or choose a unified level if specified in multiple places (definition vs options) ^{267 268}. *Error Handling*: Throws `ConfigurationError` if any required part is missing or invalid (e.g.,

missing plugin name, invalid profile structure) ^{269 270}. These exceptions are caught by `validate_settings()` and turned into a report object; if errors exist, CLI will not proceed ²²². The config loader ensures every sink has mandatory security and determinism levels, raising error if not provided ^{46 271} (this is a deliberate design to force compliance-related data to be explicit).

- **Experiment Orchestrator** (`elispeth.core.orchestrator.ExperimentOrchestrator`) – *Purpose:* High-level coordination of a single experiment's execution. *Responsibilities:* Initialize experiment context and all plugins; manage an `ExperimentRunner` that will process the data; handle security context and plugin injection ^{272 73}. It effectively bridges the input (datasource & config) with the processing (runner) and output (sinks) ². *Key Interfaces:*

- **Constructor:** Accepts `datasource`, `llm_client`, `sinks` (list), `config` (`OrchestratorConfig`), optional `rate_limiter`, `cost_tracker`, and context info (experiment name, suite_root path, config_path) ^{273 274}.
- In `__init__`, it computes a combined `security_level` for this experiment from datasource and LLM (taking the more restrictive level) ⁴⁸. It then creates a `PluginContext` object encapsulating experiment metadata (name, kind "experiment", that security level, provenance info, plus suite path) ¹²⁵. This context is passed to all plugins so they can know where they belong (for logging or for further restrictions).
- It attaches the security context to the `rate_limiter` and `cost_tracker` if they exist by calling `apply_plugin_context` on them ^{75 275}, meaning these control objects also carry the experiment identity (helpful if one rate limiter is shared across experiments to tag usage).
- It then instantiates **plugins**: iterates over each plugin definition list in `OrchestratorConfig` (row, aggregator, validation, early_stop) and calls the respective factory: `create_row_plugin(defn, parent_context=experiment_context)` etc. ^{73 276}. The `create_*_plugin` functions look up the plugin by name and call its factory with provided options, returning an instance. The orchestrator collects these in lists (or None if none defined).
- Similarly, it prepares LLM middleware: calls `create_middleware(list_of_definitions, parent_context=experiment_context)` which likely uses a registry to create each middleware instance with context ¹⁰³.
- Finally, it creates the `ExperimentRunner` (unless one was passed for testing) with all these components and settings ^{123 77}. It passes into the runner: the LLM client, sinks, prompt templates (system & user strings), prompt_fields list, criteria list, row/agg/val plugin instances, controls (`rate_limiter`, `cost_tracker`), experiment name, `retry_config` (dict or None), `checkpoint_config`, middleware list, prompt_defaults (e.g., a dict of default values for template variables), `concurrency_config` (dict or None), the resolved `security_level`, `early_stop` plugins & config ^{191 119}. This is a comprehensive injection of configuration. After constructing the runner, it sets an attribute `plugin_context` on the runner to attach the context (some plugins or controls might access `runner.plugin_context` to know their context) ²⁷⁷.
- **run() method:** This is the orchestrator's main operation. It simply loads the data via datasource and delegates to the runner:
 - Calls `df = self.datasource.load()` to retrieve input DataFrame ⁸⁵. (If this fails, exception will propagate; ideally, a DataSource should handle its own errors and maybe raise `ConfigurationError` or `IOError`, which could be caught upstream.)
 - If `max_rows` is set in config (to limit processing), it truncates the DataFrame to that number of rows ¹³³. This is helpful to avoid large runs or for debugging.
 - It extracts the prompt templates from config (system and user) into local variables (though in code it reassigns to runner later, presumably to ensure runner has latest – but

since runner was initialized with those, this step might be redundant or intended for scenario where you re-use runner across runs, which isn't typical).

- Then simply calls `payload = runner.run(df)` ¹³⁴. The runner returns the detailed results payload, which `ExperimentOrchestrator.run` then returns to the caller (CLI or suite runner).
- It doesn't itself iterate rows or handle exceptions; that's all inside `ExperimentRunner`. The orchestrator's role at runtime is just to kick it off and possibly apply any global pre/post if needed (currently none besides row limiting).
- **Dependencies:** Orchestrator depends on the plugin registry functions and on `ExperimentRunner` class ¹²¹ ¹²². It also relies on `PluginContext` and security utilities. It's tightly coupled with `ExperimentRunner` in that it needs to know all its constructor args and pass them correctly.
- **Configuration:** Orchestrator uses `OrchestratorConfig` to get all plugin definitions and settings. It doesn't read global config itself except what's handed to it via that object and explicit args.
- **Security & Context:** It is the component responsible for labeling the experiment with a security level and ensuring every plugin and sink created gets tagged with it ⁴⁸. For example, right after building a sink in suite runner, they call `apply_plugin_context` to tag sinks with context including security level, something the orchestrator likely also ensures for the global sinks list via the plugin context (though in orchestrator's code we see context for `rate_limiter`, `cost_tracker`, not explicitly for sinks – possibly sinks get context applied in suite scenario or by registry on creation).
- In summary, **ExperimentOrchestrator** is a **facade** that configures and runs one experiment. It's designed with a minimal public API (basically just `run()` after construction). It encapsulates a lot of heavy lifting in its initialization, connecting all parts so that when `run` is called, the pipeline is ready to execute. Evidence of it coordinating various parts is the series of `create_*` calls in `init` ⁷³ ⁷⁴ and the passing of context to each part.
- **Experiment Suite Runner** (`elasticsearch.core.experiments.ExperimentSuiteRunner`) –
Purpose: Manage execution of a suite of experiments (multiple experiments defined in config, possibly with baseline). *Responsibilities:* Merge hierarchical configurations (suite defaults, prompt packs, experiment overrides) for each experiment, instantiate separate `ExperimentRunner` instances for each experiment, and orchestrate their execution in sequence ⁷⁸ ²⁷⁸. It also handles baseline logic (ensuring baseline experiment runs first and computing cross-experiment comparisons) ¹⁵¹ ⁸⁰. *Key Methods:*
 - **build_runner(config: ExperimentConfig, defaults: dict, sinks: list):** This method creates an `ExperimentRunner` for a single experiment, given that experiment's specific config and the suite-level defaults ⁷⁸. It performs a three-layer merge of configuration: (1) suite defaults, (2) prompt pack defaults (if the experiment or suite references a named prompt pack), and (3) experiment-specific config (highest priority) ²⁷⁹ ²⁸⁰. A helper `ConfigMerger` is used to do field-by-field merging without manual error. For example, it merges prompt templates (system and user) such that experiment-specific ones override defaults if present ²⁸¹ ²⁸². It similarly merges lists of plugins (some config style might allow adding plugins on top of defaults, or replacing them; the code uses a special `merge_plugin_definitions` to handle that, often pack plugins + default + experiment-specific all combined) ¹²⁹. It also merges concurrency and early-stop settings and normalizes them as in single orchestrator (ensuring `early_stop_config` is turned into plugin defs if needed) ¹⁰⁰ ²⁸³.
 - It then determines the **security_level** for the experiment by taking the most restrictive of: experiment's own security setting (if any), prompt pack's (if used), and suite default's ²⁸⁴ ²⁸⁵. This parallels the logic in orchestrator (datasource vs llm) but here it's about config declarations.

- It creates an `experiment_context` (PluginContext) for this experiment with that security level and experiment name ²⁸⁶. Then it applies this context to the sinks list that will be used for this experiment (since sinks might be reused across experiments, they need to be re-contextualized for each run) ²⁸⁷. It sets each sink's `_elspeth_sink_name` and security level accordingly ²⁸⁸.
- Next, it instantiates plugin instances for this experiment (similar to orchestrator init but within the merge function): calls `create_row_plugin` etc. for each list of plugin defs it prepared ²⁸⁹ ²⁹⁰. Also, it creates LLM middleware instances via `_create_middlewarees` (which wraps `create_middleware`, caching shared ones) ⁸².
- It then handles **controls**: `rate_limiter` and `cost_tracker`. The suite can define a default `rate_limiter` config or object; the experiment can override or provide new config. The code checks if a specific config is given for this experiment (`rate_limiter_def` from merge), then either create a new instance with `create_rate_limiter` (and attach context) or reuse the one from defaults if present ²⁹¹. Similarly for `cost_tracker`. Reusing from defaults means if the suite had a global cost tracker object (like to accumulate cost across all experiments in suite), it will use that one for all experiments (after context applied) ²⁹².
- After validation (ensuring prompts are not blank after merge – raising error if system or user prompt ended up empty after all merges) ²⁹³, it constructs an `ExperimentRunner` just like orchestrator does, with the merged configuration values ⁷⁹ ²⁹⁴. It attaches `plugin_context` to it ²⁹⁵.
- The returned runner instance is fully configured for this specific experiment.
- **_instantiate_sinks(defs)**: If an experiment defines custom sinks (via `sink_defs` in config or in its prompt pack), this method creates concrete sink instances from those definitions (using `sink_registry`) ²⁹⁶ ²⁹⁷. It ensures each sink has security and determinism (or raises) and tags the sink with `_elspeth_plugin_name` and `_elspeth_sink_name` for identification ²⁹⁸ ²⁹⁹. This is used in the run loop to get an experiment-specific sink list.
- **run(df, defaults=None, sink_factory=None, preflight_info=None)**: Coordinates running all experiments in the suite on the same DataFrame ³⁰⁰. It receives the DataFrame (already loaded by CLI once, to avoid reloading for each experiment) and a `defaults` dict which usually comes from orchestrator config (or assembled in CLI `_run_suite`) ³⁰¹. The `sink_factory` parameter is interesting: CLI passes a lambda that clones base sinks for each experiment, ensuring outputs are separated (e.g., separate CSV per experiment by renaming file) ³⁰² ¹⁰⁹. If no `sink_factory` is provided, it uses the suite's global sinks by default.
 - It first builds the list of experiments to run: if a baseline experiment is defined, put that first, then all other experiments (excluding baseline if repeated) ³⁰³.
 - It prepares a `suite_metadata` list which contains metadata (like temperature, max_tokens, is_baseline flag) for each experiment in the suite ³⁰⁴. This info can be used by middleware or report generation to know context of each experiment (e.g., for logging or comparing settings).
 - It also constructs a `preflight_info` dict if not provided, including number of experiments and the baseline's name if present ³⁰⁵.
 - It keeps a dict `notified_middlewarees` to track which middleware have been globally notified of suite start (to avoid calling `on_suite_loaded` multiple times on the same shared middleware) ³⁰⁶.
 - Then it iterates through each experiment in the list:
 - Determine the effective prompt pack name and load the pack from defaults if any ³⁰⁷.
 - Determine sinks to use for this experiment: priority to experiment-specific `sink_defs` (if present in experiment config) -> else if prompt pack defines sinks -> else if suite defaults had `sink_defs` -> else use `sink_factory(experiment)` or default global sinks ³⁰⁸. This logic allows per-experiment outputs if desired, otherwise clones from base.

- Build an ExperimentRunner for this experiment by calling `self.build_runner(experiment, {**defaults, "prompt_packs":..., "prompt_pack": pack_name}, sinks)` ¹⁵⁰. It merges the defaults with prompt pack and experiment config as described.
- Retrieve the `experiment_context` from the runner (or create a fallback one if missing, but since `build_runner` sets `plugin_context`, it should be there) ³⁰⁹.
- Notify all middlewares about experiment start: For each middleware in `runner.llm_middlewares`:
 - If it has `on_suite_loaded` and hasn't been called yet for this suite, call it with `suite_metadata` and `preflight_info` ³¹⁰ (then mark it called in `notified_middlewares`).
 - If it has `on_experiment_start`, call that with experiment name and some config (like temperature, is_baseline flag) ³¹¹. This hook can be used for logging or adjusting internal state per experiment.
- Then execute the experiment: `payload = runner.run(df)` ³¹². This is similar to orchestrator doing run, but here we keep track of baseline.
- If this experiment was baseline (or is the one identified as baseline in suite), store its payload as `baseline_payload` for later comparisons ^{313 314}.
- Store the results: put `results[experiment.name] = { "payload": payload, "config": experiment }` in a results dict keyed by experiment name ³¹⁵. So the final output of `suite_run` is a dict mapping each experiment to its payload and config.
- Call `on_experiment_complete` on each middleware that has it, passing experiment name, its payload, and basic config info ¹⁰⁵. This allows middleware to log or analyze the outcome right after each experiment finishes.
- After each experiment except baseline, perform **baseline comparisons**: If we have a `baseline_payload` and `baseline_plugin_defs` (compiled from defaults+pack+experiment), instantiate each baseline plugin via `create_baseline_plugin` and call `plugin.compare(baseline_payload, payload)` ^{80 101}. Collect any diffs returned into a `comparisons` dict. If any comparisons found, attach them to payload (`payload["baseline_comparison"] = comparisons`) and also to results dict under that experiment's entry ^{101 316}. Then call `on_baseline_comparison` on middlewares if defined, so they can react to the comparison result (e.g., log a summary) ³¹⁷.
- After all experiments loop, call `on_suite_complete` on all middlewares that had `on_suite_loaded` (meaning they care about suite lifecycle) ³¹⁸.
- Return the `results` dictionary for the suite ³¹⁹.
- The SuiteRunner thus provides an aggregated view of multiple experiments' outputs, preserving each experiment's context and enabling cross-experiment analysis.
- Dependencies: It uses `create_rate_limiter`, `create_cost_tracker` from controls just like orchestrator. It uses the plugin registry for baseline plugins as well ³²⁰. It heavily uses the `ExperimentRunner` class for each experiment, demonstrating an internal dependency on that component. It also uses `ExperimentConfig` and `ExperimentSuite` data classes (likely representing parsed YAML for experiments) from `experiments.config` module ³²¹.
- Configuration: The SuiteRunner is configured with a loaded `ExperimentSuite` object (which contains multiple `ExperimentConfigs` and possibly a baseline designation) and global settings (llm client, base sinks, etc.) on initialization ³²². It then takes runtime defaults and DataFrame to run.
- Notable aspects: It **ensures deterministic behavior** by reusing LLM middleware across experiments where appropriate (to simulate how in production a single client might handle consecutive requests) ⁸². It also ensures that experiments in a suite share global controls if intended (like one RateLimiter controlling the overall rate for all experiments in that suite). This is

important: if you run 3 experiments sequentially with separate rate limiters, each might send X requests per minute; but if one global limiter is used, you won't exceed overall limits. The code as written will reuse the same `rate_limiter` object for all experiments if suite defaults provided a constructed one (which would carry state), or if not, they each get their own (which is acceptable if experiments are sequential, though if one wanted to simulate parallel experiments, it's not currently parallelizing across experiments anyway, they run sequentially).

- The baseline comparison feature is a core differentiator of suite vs single runs – it allows direct metric comparisons which are essential for A/B testing outcomes. The code's design to incorporate baseline plugins means it's easy to extend: e.g., one baseline plugin might compute statistical significance of differences (pingouin or statsmodels could be used for that).

- **Experiment Runner** (`elispeth.core.experiments.ExperimentRunner`) – *Purpose*: Execute the actual experiment logic on each data record, applying prompts and collecting outputs. *Responsibilities*: Orchestrate per-row processing (prompt rendering, LLM invocation, plugin application) and assemble the result payload for an experiment ³²³ ³²⁴. It is where concurrency and retry are implemented, making it the workhorse of runtime. *Key Interface (Dataclass fields)*:

- It's defined as a dataclass with fields for all configuration needed: `llm_client`, `sinks`, `prompt_system` (str), `prompt_template` (str), `prompt_fields` (list of columns to include), `criteria` (list of criterion definitions), lists of plugin instances (`row_plugins`, `aggregator_plugins`, `validation_plugins`), optional `rate_limiter` and `cost_tracker`, `experiment_name`, optional `retry_config` and `checkpoint_config`, plus internal fields for compiled templates and tracking (like `_compiled_system_prompt`, `_compiled_user_prompt`, `_compiled_criteria_prompts`) ²³² ³²⁵. Also `llm_middlewares`, `concurrency_config`, and fields for security and determinism level are stored ³²⁶. Early stop plugins and config are stored, and internal sync primitives `_early_stop_event` and `_early_stop_lock` are created as needed ³²⁷.

- **run(df)**: This is the primary method to run the experiment on the given DataFrame ³²³. Steps (some already described in data flow):
 - Initialize early stop mechanism via `_init_early_stop()`: set up `_active_early_stop_plugins` list and Event/Lock if needed (create plugin instances for early_stop if config given but no plugin list) ³²⁸ ³²⁹. It resets any previous early_stop state so a fresh run starts clean ³³⁰.
 - Initialize checkpoint tracking: if `checkpoint_config` is set, load existing checkpoint file to get `processed_ids` (set of IDs) and determine `checkpoint_field` (the column name to use as unique ID, default "APPID") ³³¹ ¹³. Also open a path to append new IDs for any newly processed rows.
 - Compile prompts: uses `PromptEngine` (or existing `prompt_engine` if injected for testing) to compile the system prompt and user prompt templates into `PromptTemplate` objects ¹⁹⁴ ³³². Also compile each criterion's template, merging any criterion-specific defaults with global defaults ¹⁶⁴. The compiled templates are stored in `_compiled_system_prompt`, `_compiled_user_prompt`, `_compiled_criteria_prompts` for reuse during the run.
 - If datasource provided a schema (df.attrs['schema']), run `_validate_plugin_schemas(schema)` to ensure the plugins can handle the data schema ⁹⁶. This likely uses jsonschema or internal logic to verify plugin compatibility (the code calls `validate_schema_compatibility` in imports ¹⁶¹).
 - Prepare to iterate: Initialize list for results and failures, and prepare list `rows_to_process`:

- Iterate over `df.iterrows()`, for each row create context dict using `prepare_prompt_context(row, include_fields=self.prompt_fields)` ¹⁶³. This likely converts row Series to a normal dict (if `prompt_fields` is None, include all columns; if a list, include only those keys).
- Determine `row_id`: if `checkpoint_field` is set and present in context, use it; else None ¹⁸³.
- If `processed_ids` (from checkpoint) is not None and `row_id` is in it, skip this row (already processed in a previous run) ¹⁸⁴.
- If `early_stop_event` is set (some prior row triggered stop), break out of loop ³³³.
- Otherwise, append `(index, row_series, context, row_id)` to `rows_to_process` list.
- Now have a list of rows to actually process (not skipped by checkpoint or early stop). Determine concurrency: call `_should_run_parallel(concurrency_cfg, backlog_size)` to decide if we run parallel or sequential ¹³⁷. By default, `concurrency_config` is None or disabled, unless user config enabled it. This function checks `config.get("enabled")` and if `max_workers > 1` and backlog above threshold ⁷. If returns True:
 - Call `_run_parallel(rows_to_process, engine, system_template, user_template, criteria_templates, row_plugins, handle_success, handle_failure, config)` ³³⁴. This spins up threads:
 - Create `ThreadPoolExecutor` with `max_workers` (from config, default 4) ³³⁵ ³³⁶.
 - Define `worker(data_tuple)` function which processes one row: calls `_process_single_row()` to get record or failure, then acquires a lock to append to results or failures list safely ²³⁷ ³³⁷.
 - The main thread submits each item in `rows_to_process` to the executor, but before submitting each, it checks rate limiter utilization: if `utilization >= pause_threshold`, sleep for `pause_interval` and re-check ³³⁸. This dynamic pacing prevents flooding if the system is near capacity.
 - Also, if `early_stop` event is set while scheduling (in case it triggers during parallel execution), break out of scheduling loop ²⁴².
 - Threads themselves also check `early_stop_event` at start of worker and simply return without processing if set after scheduling but before execution started on that thread ³³⁹.
 - The `handle_success` and `handle_failure` closures essentially do the same as in sequential loop: append results and manage checkpoint. In code, they directly append to shared lists under lock and call `_maybe_trigger_early_stop` for each success.
 - `_run_parallel` doesn't explicitly join threads because the `with ThreadPoolExecutor` block ensures all submitted tasks complete before exiting ³⁴⁰ ¹⁵⁴. Once done, it returns to proceed with aggregation.
 - If concurrency is not enabled or backlog is below threshold:
 - Use a simple loop: for each `(idx, row, context, row_id)` in `rows_to_process`:
 - Check `early_stop` event at top; break if set ³⁴¹.
 - Call `_process_single_row(engine, system_template, user_template, criteria_templates, row_plugins, context, row, row_id)` synchronously ³⁴².
 - If a `record` comes back (not None), call `handle_success(idx, record, row_id)`; if a failure dict comes, call `handle_failure(failure)` ³⁴³. These do same as above: append to records list and checkpoint, or append to failures.

- After parallel or sequential loop, sort `records_with_index` by index and extract the sorted records (to maintain original input order in output) ²⁰⁷ .
- Assemble `payload`: put `"results": results` (list of result dicts). If failures list non-empty, include `"failures": failures` ³²⁴ .
- Compute aggregates: iterate aggregator plugins and call `plugin.finalize(results)`, if a plugin returns a derived value, add it to an `aggregates` dict (keyed by plugin name) ⁹² . If any aggregates, attach `"aggregates": aggregates` to payload ³⁴⁴ .
- Prepare `metadata`: a dict with `"rows"` and `"row_count"` equal to number of results (for convenience) ²⁰⁸ . Compute `retry_summary`: iterate over all results and failures, summing up total requests, total retries, and count how many were exhausted (failures) ²¹ . If any retry info found, include `"retry_summary"` in metadata ³⁴⁵ . Include aggregates in metadata if present ³⁴⁶ . If a `cost_tracker` is used, get its summary (like total tokens & cost) and include as `"cost_summary"` in both payload and metadata ³⁴⁷ . If any failures occurred, include them in metadata as well (`metadata["failures"] = failures`) ³⁴⁸ .
- Determine the effective security and determinism levels for the data: If DataFrame had attributes for those (or none), combine with the experiment's configured level by `resolve_security_level(self.security_level, df_security_level)` ¹⁴⁸ and similarly for determinism ³⁴⁹ . This sets `_active_security_level` and `_active_determinism_level` on the runner and adds them to metadata. This ensures if input data had a classification higher than the experiment classification, the output is labeled with the higher (most restrictive) classification.
- If `early_stop` triggered (`_early_stop_reason` is set), add `"early_stop": reason` to both payload and metadata ³⁵⁰ .
- Attach the metadata dict to payload as `payload["metadata"] = metadata` ³⁵¹ .
- Instantiate the `ArtifactPipeline` with `self._build_sink_bindings()` ¹¹⁷ . `_build_sink_bindings` creates a `SinkBinding` for each sink in `self.sinks`, capturing its id, plugin name, instance, artifact config, original index, and `security_level` ^{352 353} . Security level is normalized for each sink (if sink had `_elspeth_security_level` attribute, use that) ¹¹⁵ . This method returns a list of `SinkBindings`.
- Call `pipeline.execute(payload, metadata)` . This triggers all sinks to run in proper order as detailed earlier. It effectively delivers the payload (and artifact store) to each sink's `write()` . For example, a `CsvResultSink` will receive the payload and possibly filter out "results" to write them.
- After pipeline execution, set `_active_security_level = None` (perhaps as a cleanup to avoid accidental reuse or to signal that pipeline finished) ³⁵⁴ .
- Return the payload dict as final output of runner ³⁵⁵ .
- **_process_single_row(...)**: Handles the actual processing of one row. Already stepped through: it renders prompts (calls `_render_prompts` to get system and user prompt strings given context) ^{356 357} , then calls `_collect_responses(rendered_system_prompt, base_user_prompt, criteria_templates, context, row, row_id)` which calls `_execute_llm` for primary prompt and each criteria (if any) and composes a record with responses ^{358 359} . `_collect_responses` uses either direct call or a separate method for multiple criteria to gather all LLM outputs. The primary response is returned as `response` , and possibly a `record` containing "responses" for criteria and "response" for main. `_collect_criteria_responses` (not fully shown) likely loops through criteria and calls `_execute_llm` for each one, building a dict of responses.
 - Then `_populate_prompt_metadata` attaches the original rendered system and user prompt (and maybe criterion prompts) to the record for traceability ¹⁴⁰ . This ensures if

needed, one can see exactly what prompt was sent for each row (helpful in debugging or audit).

- `_attach_retry_metadata(record, primary_response)` pulls out the retry info from `primary_response` (which `_execute_llm` puts under `response["retry"]`) and sets `record["retry"] = that info`, and also something like `record["response"]["metrics"]["attempts_used"] = attempt_count` in the last attempt, to record how many attempts it took ¹⁶⁸.
 - `_apply_row_plugins(record, row_plugins)` iterates each row plugin instance, calls `plugin.process_row(record["row"], record.get("responses") or record.get("response"))`, merging its return into `record["metrics"]` if not None ⁸⁸ ⁸⁹. Exceptions in plugins here bubble up to broad catch in run loop (so one failing row plugin would mark row as failure).
 - `_apply_security_level(record)` likely tags the record with current active security level (taking into account any plugin changes? Possibly not needed since security set at experiment level).
 - Returns (record, None) on success, or (None, failure_dict) if something went wrong. It has multiple except blocks:
 - Catches `PromptRenderingError` and `PromptValidationError` specifically and returns a failure dict with the context row and error message ³⁶⁰.
 - A broad `except Exception` to catch any other error; it creates a failure with context, error string, timestamp ³⁶¹ ³⁶² (and possibly attaches type or history if present) before returning it.
- **`_execute_llm(system_prompt, user_prompt, metadata)`:** This is perhaps the most intricate piece. It handles the retry loop:
- Initialize `max_attempts` (from `retry_config` or default 1), `delay` (initial backoff), `backoff` factor, and attempt counter ³⁶³.
 - Loop from attempt 1 to `max_attempts`:
 - Mark attempt start time.
 - Construct an `LLMRequest` object with `system_prompt`, `user_prompt`, and `metadata` (which includes "attempt": attempt number) ³⁶⁴. Save `last_request`.
 - Apply each `before_request` middleware in order to the request, possibly transforming it or replacing it ¹⁰⁶.
 - If a `rate_limiter` is present, call `acquire_context = self.rate_limiter.acquire({...})` providing some metadata like experiment name and perhaps attempt. If it returns a context (like a context manager controlling a token bucket), enter it (with `acquire_context`) around the call to ensure the call counts toward the rate limit usage ³⁶⁵.
 - Call the LLM client's `generate(system_prompt=request.system_prompt, user_prompt=request.user_prompt, metadata=request.metadata)` to get a response. This is the actual API call (could raise exceptions on network or API error) ³⁶⁶.
 - Apply `after_response` middlewares in reverse order on the response (maybe filtering or adding to it) ⁵⁷.
 - If a `cost_tracker` is present, call `cost_tracker.record(response, request.metadata)` which returns any cost metrics (like tokens, cost). If metrics returned, merge them into `response["metrics"]` ¹⁶⁷.
 - Run validations on the response via `_run_validations(response, request, row_context)` ³⁶⁷. If a validation plugin raises an exception, it will be caught by broad catch below, unless `_run_validations` handles them internally (not clear, maybe it raises a `PromptValidationError` which would be caught by outer try).

- On success, create `attempt_record` dict with attempt number, status "success", and duration time ³⁶⁸. Append it to `attempt_history` list. Set `response["metrics"]` `["attempts_used"] = attempt` (so if it took 3 attempts, `attempts_used=3` in metrics) ³⁶⁹.
- Also set `response["retry"] = { "attempts": attempt, "max_attempts": max_attempts, "history": attempt_history }` for full trace ³⁷⁰.
- If a `rate_limiter` is present, call `rate_limiter.update_usage(response, request.metadata)` to record the outcome (some limiters might count tokens used or note that a call happened, etc.) ³⁷¹.
- Return the `response` dict to caller (meaning success, break out of loop) ³⁷².
- If an exception occurs in the try (covering everything from `before_request` to validations):
 - Catch it as `exc`. Save `last_error = exc`.
 - Create an `attempt_record` for this attempt with status "error", duration, and error message and type ³⁷³. If not final attempt, also calculate `sleep_for = delay` (the current delay, 0 if none set) and put that in `attempt_record` as "next_delay" ³⁷⁴. Append to `attempt_history`.
 - If `attempt >= max_attempts`, break out of loop (exhausted all tries) ³⁷⁵; else:
 - If sleep delay > 0, sleep that duration ³⁷⁶.
 - If backoff factor > 0, multiply delay by backoff for next attempt (exponential backoff) ³⁷⁷.
 - Continue loop to next attempt.
- After loop, if `last_error` is not None (meaning all attempts failed):
- If `last_request` is not None, attach `_elspeth_retry_history` and `_elspeth_retry_attempts/max_attempts` attributes to the exception object for later inspection ¹⁹⁹.
- Call `_notify_retry_exhausted(last_request, last_error, attempt_history)` (not shown, likely informs middlewares or logging that we gave up).
- The code snippet doesn't show a `return` or `raise` explicitly after the loop for failure, but likely it will raise `last_error` so that it gets caught by the outer `_process_single_row` and treated as failure.
- Dependencies: uses `LLMRequest` and `LLMMiddleware` from `core.protocols`, which defines interface for middleware (with `before_request/after_response` methods) ³⁷⁸. The LLM client is expected to implement `generate` as per `LLMClientProtocol` (some plugins might call external API or return a dict directly, as dummy did) ⁸⁶.
- The design ensures that if a final attempt fails, it doesn't return a response but instead triggers failure handling in the caller; if any attempt succeeds, it returns immediately with success and includes all attempts in history (for auditing partial failures if any).
- Rate Limiter usage pattern: If defined and has a utilization method, they throttle the scheduling of threads; also the actual request acquisition uses context management. The actual `RateLimiter` implementation in `core.controls.rate_limit` likely uses a token bucket or counting semaphore pattern. The code in runner coordinates with it carefully. Similarly, `cost_tracker`'s update after getting response ensures cost per response is accounted and can be summarized.

Summarizing `ExperimentRunner`: it's a self-contained execution engine that, once configured with all needed objects, can run an entire experiment and yield full results. It's perhaps the most critical component for correctness and performance. It deals with threads, external calls, error handling, and

ensures that at the end, outputs are correctly passed to sinks. The design isolates this complexity so that orchestrator and suite logic can be simpler.

- **Plugins** – We have many plugin types; listing major ones:

- **DataSource Plugins:**

- **CSVLocalDataSource** – likely reads a local CSV file path using pandas. *Interface:* It implements `load()` which returns a DataFrame ⁸⁴. *Configuration:* Options include `path`, possibly `encoding`, etc. *Dependencies:* Uses `pandas.read_csv`. *Key features:* Could attach `df.attrs['security_level']` if config specifies, e.g., if the config for datasource had `security_level: SECRET`, the DataFrame can carry that. It might also attach a `schema` by inferring from columns. For example, maybe reading a CSV with a header and sample data, it could construct a simple schema describing column names and types. In absence of explicit, it may leave schema None.
- **AzureBlobDataSource** – likely uses azure-storage-blob to fetch a blob (CSV or Excel) from Azure. *Interface:* same `load()`. *Configuration:* Options might include container, blob name, or perhaps a connection string or use of azure-identity for auth. *Function:* It would call `BlobClient.download_blob()` and then pass to pandas (if CSV) or if it has a `format` option. The docs mention that initial CLI was focusing on Azure Blob ingestion ⁵⁸, likely via such a plugin. *It almost certainly attaches a security label:* Many orgs classify data at rest, so config might specify e.g. `security_level: TopSecret` for that blob source, which the plugin would propagate to DataFrame attributes. The dependency analysis notes "monitor credential escalation CVEs in azure-identity" and TLS enforcement ⁵⁵, implying blob plugin uses azure-identity to get credentials (which might use environment credentials or managed identity, thereby no secrets in config).
- **MockDataSource** – they might have a dummy for testing that generates a small DataFrame from some static data (like the test's DummyDatasource does inline) ⁸⁴. Possibly not needed as plugin because tests just define a dummy class.
- *General behavior:* Data sources should ideally handle **data validation** (ensuring required columns are present, or maybe applying any pre-processing, e.g., type conversions). The code references `SchemaViolation` and an `on_schema_violation` setting (which in runner is "abort" by default) ³⁷⁹. Possibly, a DataSource might set up `df.attrs['schema']` with something like a dictionary of expected fields and types, or an instance of a schema class. Then, a validation plugin or the runner itself will check if row plugins or LLM prompt fields align with this schema. They also might filter out rows with missing required data (or mark them for skip) – but likely that's left to validation plugins, not data source.

- **LLM Client Plugins:**

- **OpenAIClient** (`plugins/llms/openai_http.py`) – integrates with OpenAI's API. *Configuration:* probably expects an API key (if not using azure-identity) and optional parameters like model name, temperature, etc. Possibly the config might allow specifying some default generation parameters (like model "gpt-3.5-turbo", temperature 0.7).
- *Dependencies:* uses the `openai` Python library (OpenAI SDK) ³⁸⁰, or possibly just `requests` if doing HTTP manually. The dependency file shows `openai>=1.12`, implying they use the official SDK which handles retries and JSON etc. The plugin's `generate()` function likely calls `openai.Completion.create()` or `openai.ChatCompletion.create()` depending on how they structure prompts (the presence of system prompt and user prompt suggests ChatCompletion context).

- *Response*: likely normalizes the OpenAI API response (which is a complex object with choices, usage, etc.) into the simpler dict format the runner expects: e.g., `{"content": "<text>", "metrics": {"tokens_prompt": x, "tokens_completion": y, "tokens_total": z}}`. Possibly also attaches model or request id as needed. The cost tracker might then convert token counts to cost if it knows pricing for that model (cost tracker could have model pricing table).
- The plugin should be careful with exceptions: OpenAI library raises `openai.error.OpenAIError` for issues (like rate limit, invalid request). Those will be caught by ExperimentRunner's try and treated as needing retry if configured. They may also raise on things like network issues (which come as `requests.exceptions` if not handled by SDK).
- **AzureOpenAIClient** (`plugins/llms/azure_openai.py`) – for Azure's variant (which uses different endpoints and possibly azure identity auth). *Config*: includes endpoint URL, deployment name of model, etc. *Implementation*: possibly uses the OpenAI SDK but configured with `api_type="azure"` and `api_base`, or uses `azure.core` libraries. The dependency analysis lumps openai lib with Azure usage, implying they might still use openai lib by setting environment variables for Azure endpoint, or use azureml-core in some way. Actually, azureml-core might be used for telemetry, not for calling the model.
- It might accept an `azure_credential` (like azure-identity's DefaultAzureCredential), especially since azure-identity is a dependency ³⁹. If so, the plugin could retrieve a token and call the Azure REST endpoint with requests (since the OpenAI SDK doesn't handle Azure AD tokens directly, I think).
- This plugin should output similar format. Possibly, it might include content safety scores if Azure returns them (if not using a separate call).
- **MockLLMClient** (`plugins/llms/mock.py`) – a dummy that returns a fixed or echo response for testing and offline mode. *Behavior*: likely simply returns `{"content": <some static content or echo of user prompt>}` quickly. The test DummyLLM did return `{"prompt": user_prompt, "meta": metadata}` or in other test `{"content": user_prompt, "metrics": {"score": 0.5}}` ⁸⁶ ³⁸¹. The real Mock plugin might be configurable to either echo input or use a small local model (like random choice from predefined outputs).
- This is crucial for offline, and code does mention using the mock when no internet, to avoid hang ⁶⁵.
- Possibly others: e.g., **AzureContentSafety** might be integrated as an LLM middleware rather than client, scanning content post-hoc. But mention is that requests lib is used by "Azure Content Safety" and "middleware" ⁴⁵, likely implemented as a middleware that calls Azure's content filtering API in `after_response`, potentially dropping or flagging harmful content.
- Each LLM client plugin must abide by `LLMClientProtocol` with a `generate()` method. They register themselves with `llm_registry` such that config `"llm": {"plugin": "openai_http", "options": {...}}` is resolved to an instance.

• Experiment (Row/Aggregator/Baseline/EarlyStop/Validation) Plugins:

- **Metrics Plugins** – Possibly they have built-in ones like:
- **ScoreExtractor** (row plugin): might extract a numeric score from LLM response if the prompt asks for a rating. In test, they allude to metrics plugin names like "score_extractor", "score_stats", etc., which CLI's `--disable-metrics` filters out ³⁸². `ScoreExtractor` (row-level) might parse the LLM content if it contains a number and return `{"score": <float>}`.

- **ScoreStats** (aggregator plugin): compute mean/median of score across all rows, output summary stats.
- **ScoreRecommendation** (aggregator plugin): maybe suggests if one model is better than another based on scores (if baseline).
- **score_delta** (baseline plugin): compare experiment's scores to baseline's and produce differences. The CLI's `_strip_metrics_plugins` specifically filters plugins by name: it knows certain plugin names to remove when metrics are disabled ³⁸³ ³⁸⁴. Names mentioned:
 - row plugin named "score_extractor"
 - aggregator plugins "score_stats", "score_recommendation"
 - baseline plugin "score_delta". So yes, those exist in `plugins/experiments/metrics.py` perhaps. *These plugins* add quantitative evaluation. Implementation: likely straightforward (score_extractor uses a regex or JSON to find a number; score_stats does `np.mean` etc using numpy or SciPy).
- **Prompt Variant Plugin** – Perhaps the config allows generating prompt variants by plugin? Actually, there is `plugins/experiments/prompt_variants.py` (seen in search results). Might be a plugin that, after baseline run, creates additional prompt variations for new experiments dynamically (just speculation).
- **Validation Plugin** – Could be two types:
 - *Content validation*: e.g., check that the LLM's content meets certain format. Possibly an example is "JSONValidator" that ensures the response is valid JSON. If not, it might raise `PromptValidationError` (caught and marks failure) ³⁶⁰. Or "SchemaValidator" if the LLM output is supposed to have certain fields.
 - *Input validation* (less likely, input is static pre-provided). They provide a `validation_plugins` mechanism parallel to row and aggregator, and also a separate notion of schema validation outside plugin context (the runner does initial check with `validate_schema_compatibility`). Possibly, more advanced validation like "no PII present in response" could be a validation plugin calling some detection logic.
- **Early Stop Plugin** – They have one named "threshold" as referenced in code ³⁸⁵. Possibly it monitors something like if a certain metric in record passes a threshold, or if certain count of failures happen. It implements `check(record, metadata=None) -> reason or None`. E.g., threshold plugin might be configured like: stop if `metric_score < 0.2` for some record (meaning performance is too bad, abort further trials), or stop after `N` rows processed. Another early_stop might be "cost_limit": config could say cost limit \$X, plugin check accumulative cost via `cost_tracker` and if exceeded returns reason. Actually, `cost_tracker` integrated within can trigger early stop indirectly: the code doesn't show `cost_tracker` hooking into early stop, but one could imagine an early_stop plugin reading cost from `cost_tracker` each row (maybe using metadata from `cost_tracker` updated usage).
- Implementation likely uses the `metadata` passed to `check` to see overall or uses known fields in record. They wrap all early_stop defs to uniform format via `normalize_early_stop_definitions` (the code does in `plugin_registry` or config) ³⁸⁶.
- The test references threshold plugin in `_init_early_stop`: they create a default def `{"name": "threshold", "options": {...}}` if `early_stop_config` is given without plugin defs ³⁸⁵. So threshold plugin likely can be configured solely by `early_stop_config` fields.
- **Baseline Plugins** – Provide cross-experiment comparisons. They are executed after all experiments by `suite_runner`. A default one `BaselineComparisonPlugin` might simply compute differences in certain aggregate metrics. For instance, if experiments produce an average score metric, baseline plugin could output the delta and maybe significance.

The code variable `comp_defs` collects `baseline_plugin_defs` from `defaults`, `pack`, and `experiment` ⁸⁰. There's mention of plugin name in baseline context in code: `plugin = create_baseline_plugin(defn, ...) diff = plugin.compare(baseline_payload, payload)` ³²⁰.

- Possibly, "score_delta" plugin is such a baseline plugin (compares baseline vs variant scores).
- If using stats, maybe a plugin uses Pingouin (which is listed as dependency for stats-agreement) to run tests like Wilcoxon signed-rank or correlation, outputting a p-value or effect size.
- **Observer/Telemetry Plugins** – They did not explicitly call these "plugins", but the LLM middleware serve that role (observing requests/responses). Also, the `cost_tracker` and `rate_limiter` are not exactly plugins, but parts of controls registry. Those are created via config but they don't have the same interface as others (though they do have `.update_usage`, etc.).

• Output Sink Plugins:

- **CsvResultSink** – Writes results to a CSV file. *Responsibilities*: Flatten each result record into a row and write to CSV. *Config options*: file path, whether to overwrite or append, and `sanitize_formulas` (bool) and `sanitize_guard` (maybe a string pattern) as seen in CLI clone function ¹¹⁰. Sanitization means: if a cell value begins with '=' or '-' or '+' (Excel formula triggers), prefix with a single quote to neutralize it (common mitigation for CSV injection). Also maybe handle special Unicode that Excel might interpret. The sink likely uses Pandas or csv module to write out. Because the CLI final output also writes CSV if asked, the sink might normally operate in dry-run during CLI head preview (i.e., not actually writing if CLI has that logic, but if user passes `--live-outputs` or specifically configures a CSV sink, it will truly write).
- The sink has a `write(results, metadata)` that probably does: if `sanitize_formulas` true, apply regex to any values matching `'^[=+\\-@]'` and prepend `'` (the CLI code suggests this, as it passes `sanitize_formulas` and a guard param to clone) ¹¹¹.
- *Security*: It requires a `security_level` and `determinism_level` in config, which pipeline checks. If sink is for Official data and record is Official, allowed. If record were Secret, this sink would have needed clearance Secret or pipeline aborts.
- **ExcelSink** – Writes results to an Excel file (openpyxl used). It might write multiple sheets: e.g., one sheet with raw results, maybe another with aggregates or charts if any. The dependency openpyxl is included in extras ³⁸⁷. *Options*: file path, perhaps template workbook to use or something. It likely uses openpyxl to create a workbook and style it. They mention in dependency analysis: "include openpyxl only when spreadsheets necessary" ³⁸⁸, and indeed they keep it in a separate extra for sinks-excel.
- Excel formulas sanitization is trickier: openpyxl by default will keep formulas formula (which could be malicious if the user opens it? Actually, Excel formulas in an .xlsx can still do things like DDE exploits, but those are more complex. The sanitization probably focused on CSV specifically because that can execute as soon as opened or in some contexts).
- **VisualReportSink** – Generates visual charts. Possibly uses matplotlib & seaborn (the extras analytics-visual includes them ⁶⁰). *Function*: Summarize results in a bar chart or confusion matrix, etc., and either save as PNG or HTML with embedded images. If HTML, there is risk of XSS if content included, but presumably it's mostly graphical or static text so low risk.

- It might produce multiple artifact files (PNG images). The artifact pipeline, if multiple images, the sink might declare produces for each image and pipeline will store them individually.
- **LocalBundleSink** – Perhaps collects all outputs into a local directory structure for archiving. Maybe moves files to a folder or zips them.
- **ZipBundleSink** – Takes all artifacts produced by earlier sinks (CSV, Excel, images) and zips them into one archive file. *Config*: output zip file path. *Consumes*: likely `ArtifactRequest` tokens like 'csv', 'excel', 'png', etc. *Produces*: one artifact of type 'zip'. It might also optionally delete the individual files after zipping to keep only archive, but likely not by default.
- **SignedSink** – Cryptographically signs either individual artifacts or a manifest of artifacts. Possibly it:
 - Collects all artifacts by consuming 'all' with mode 'all' (meaning get every artifact from store) ³ ³⁸⁹.
 - Produces maybe a detached signature file (like .sig or .asc) or a signed manifest (like a JSON file listing each artifact's hash and a signature block). The `core.security.signing` likely provides a function to sign data given a private key. Perhaps config includes a path to a private key or uses an HSM or Azure Key Vault via azure-identity (less likely unless integrated).
 - The signature ensures integrity: e.g., an auditor can verify that outputs were not tampered with after run. If license or classification is important, maybe they sign to prove classification labeling was in place at time of run.
 - If they didn't include a crypto library, perhaps they rely on GPG via subprocess or not implemented fully yet. There's mention in docs of signing but we didn't see a direct crypto dependency. Possibly they use Python's `hashlib` and store a hash (not a true signature).
- **RepositorySink** (GitHub/Azure DevOps) – Interacts with remote repos:
- **GitHubRepoSink**: *Function*: for each artifact (like a CSV or Excel file), create or update a file in a GitHub repository (maybe via GitHub API, using `requests`). Config would include repo URL, branch, maybe a GitHub token (they likely expect to find token in environment or config, e.g., GITHUB_TOKEN).
 - Dry-run Mode (default) prevents actual push – it might just log what it would do or write to a temp local path. The CLI toggles dry_run off only if `--live-outputs` given ³⁴.
 - On `write`, if dry_run is true, it possibly does nothing or prints a message. If false, it uses `requests` to call GitHub contents API to push new file or update if exists (the commit will have an audit trail).
- **AzureDevOpsRepoSink**: similarly uses Azure DevOps Git API or maybe creates a pull request in Azure DevOps.
- These sinks are potentially sensitive because pushing data outside the immediate environment requires careful security level checks. The pipeline ensures classification is allowed ³. E.g., if data is Official-Sensitive and your GitHub repo is not approved for that level, one would not allow this sink's clearance to be set accordingly. The config must specify `security_level` for the sink and the system should only allow if it's equal or above the data's level.
- They also likely allow specifying a path or file naming convention in the repo. The CLI clones base sinks by injecting experiment name into filenames for suite runs ¹⁰⁹ (for CSV sink example, base path "results.csv" becomes "experimentA_results.csv"). Possibly repository sink might commit into a folder per experiment or branch per experiment to keep things organized.

- **AnalyticsReportSink** (perhaps the SuiteReportGenerator uses similar logic, but maybe they also have a sink that directly produces a combined report like a Markdown or HTML summary).
- Sinks have to abide by `ResultSink` protocol with `write(results, metadata, artifacts?)`. We see pipeline giving both payload and metadata to `pipeline.execute`, but how it passes to sink is internal. Possibly, pipeline calls `sink.write(filtered_artifacts, metadata=metadata)` where `filtered_artifacts` is a dictionary of artifacts it consumes resolved from ArtifactStore (the `ArtifactPipeline.resolve_requests` does produce a dict of token->list of Artifacts) ³⁹⁰ ³⁹¹. So maybe `sink.write(artifact_map, metadata)` rather than full results. Alternatively, simpler: each sink is given the entire payload and it internally chooses relevant parts (like in CSV sink's implementation, it would look at `payload["results"]`). However, the presence of artifact config suggests a more decoupled approach.
- Each sink must declare what it produces (for store indexing). For example, CSV sink might produce type 'csv' with a name like the file name (so artifact id can be used if needed), Excel sink type 'excel', repository sink might produce nothing (or produce a commit artifact? They likely don't store commit ID as artifact, but they could).
- Sinks run in the same process after experiment finishes, sequentially as pipeline orders them. If a sink fails (e.g., network error pushing to GitHub), pipeline stops and raises error. At that point, the orchestrator or suite runner would not catch it explicitly, so it would propagate to CLI causing the CLI to terminate with error, possibly after other sinks had done their job. There's likely guidance that in such case, user should examine logs and maybe retry outputs (there's no built-in retry for sinks except possibly inherent in requests' logic or if one wraps sink in retry plugin which is not indicated).

• Control Components:

- **RateLimiter** (e.g., `TokenBucketRateLimiter`): *Purpose*: enforce call rate limits to external LLMs. *Implementation*: likely token bucket or leaky bucket: config can have `rate: N calls per interval`. The `acquire(metadata)` returns a context manager that waits if necessary until a token is available, then yields (and returns the token when context exits). The `utilization()` returns fraction of capacity currently used (for pause logic) ⁸. After each call, `update_usage(response, metadata)` might remove a token or update count.
- Perhaps `azure_openai` plugin uses Azure's rate-limit headers to calibrate rate limiter usage, but not sure if that's integrated.
- Config might allow specifying a specific QPS or TPM and maybe burst size.
- It's not a plugin in `plugin_registry` but built via `create_rate_limiter` function (in `core.controls.registry` or similar) which picks an implementation (like if config says type: "token_bucket" with options).
- It's applied globally or per experiment depending on how configured. E.g., one might set one limiter for the suite to throttle total calls across all experiments (the code allowed passing a base rate_limiter to all experiments by referencing defaults).
- **CostTracker**: *Purpose*: accumulate costs (monetary or token usage) from LLM calls. Possibly it sums `response["usage"]["total_tokens"]` or uses known model pricing to accumulate a cost total. The `cost_tracker.summary()` returns a summary (like `{"total_tokens": X, "estimated_cost_usd": Y}`) which runner adds to payload ³⁴⁷. It might also break down cost by model if multiple used.

- Implementation: For OpenAI, each response's `usage` field has tokens, then the cost per 1K tokens for model is known. The dependency analysis mentions hooking up cost tracking and merging metrics ¹⁶⁷.
- Could also track compute time cost, but likely focusing on API cost.
- Config could allow specifying budgets (like stop if cost > \$10, though `early_stop` plugin might handle that by checking `cost_tracker`).
- Probably a simple class with methods `record(response, metadata)` and `summary()`. They attach it with context like others for possible auditing (maybe `cost_tracker` has a `.name` or `.security_level` though not necessarily needed).
- **Telemetry Middleware:** They mention Azure ML telemetry in dependency extras (`azureml-core`) ⁵⁴. Possibly there's a middleware `AzureMLTelemetryMiddleware` (`middleware_azure.py`) that logs prompt & response metadata to Azure ML (e.g., an `MLRun` or `Experiment` in Azure Machine Learning for tracking). `azureml-core` is large but they included it as optional, meaning that if one wants to integrate with Azure ML (perhaps to log metrics or artifacts to an Azure ML workspace experiment), this middleware can do that on `on_experiment_complete` for example. Or maybe it uses `azureml` to gather telemetry about usage of model endpoints.
- If not used, `azureml-core` not installed. If used, configuration might specify workspace details or assume environment already set (like running inside an Azure ML run context, the middleware might just log metrics via `run.log`).
- This is beyond core function but included for enterprise observability (the README mentions enterprise observability – likely referencing hooking into existing monitoring systems) ³⁹².

Each of these components (plugins, sinks, controls) is designed to be **replaceable and extensible**. E.g., an organization could implement a custom `DataSource` (maybe pulling from a database) and register it via `datasource_registry.register("my_database", MyDatabaseSource)`. The code supports that pattern (in tests they register a plugin to simulate new ones) ¹⁴.

Finally, listing the major components in summary form: - **CLI & Config:** `cli.py` and `config.py` - user interface and configuration builder. - **Core orchestrators:** `ExperimentOrchestrator`, `ExperimentSuiteRunner` - coordinate experiments (single or multiple). - **Execution engine:** `ExperimentRunner` - runs an experiment end-to-end with concurrency, controls, etc. (the **heart** of processing). - **Plugin system:** - Data sources (`core.datasource_registry` and implementations in `plugins/datasources/*`), - LLM clients & middleware (`core.llm_registry` and `plugins/llms/*`, `plugins/llms/middleware*.py`), - Experiment plugins (`row`, `agg`, etc., in `core.experiments.plugin_registry` and `plugins/experiments/*`), - Sinks (`core.sink_registry` and `plugins/outputs/*`). Each registry maps plugin names to factory functions. Factories create instances given normalized options (with security level, determinism included) ³⁹³ ¹²⁴. - **Security & Compliance:** - `core.security.signing` - handles cryptographic signing of artifacts (if implemented). - `core.security.resolve_security_level` & `normalize_security_level` - logic to ensure consistent labeling (for example, accept synonyms or different case for levels and standardize them) ³⁹⁴ ³⁸⁹. - The classification scheme ("OFFICIAL", etc.) is known to the system; likely define allowed levels and ordering somewhere (the pipeline uses `is_security_level_allowed(producer_level, consumer_level)` ³⁹⁵, which probably knows an order like `Official < Secret < TopSecret` or similar and returns `True` if `consumer_level` is `>= producer_level`). - `ConfigurationSecurity`: config might allow default security level for entire run or separate ones for different sections (like each plugin can be given a security level that must be `>=` overall? Not sure if they allow that granular, but at least sinks require explicit levels). - Logging and auditing hooks, e.g., `audit_logging.md` suggests maybe there's an audit log plugin or simply that logs themselves are considered audit trail (with `level=INFO` logs capturing key events). - **Testing**

components: Not deployed but part of component list: e.g., `DummyDatasource`, `DummySink` used in tests verify that orchestrator and runner integrate properly ¹⁵³ ³⁹⁶. This indicates the architecture is testable by injecting dummy components (helpful for isolation and verifying that e.g., row plugins are called, sink gets correct metadata, etc.).

Each component interacts through well-defined interfaces, mostly Python abstract classes or protocols (PEP 544). For example, in `core.protocols` they define Protocol classes:

```
class DataSource(Protocol):
    def load(self) -> pd.DataFrame: ...
class ResultSink(Protocol):
    def write(self, results: Any, *, metadata: dict[str, Any] | None = None)
    -> None: ...
```

These protocols ensure components can be treated uniformly by orchestrator/runner. This fosters extensibility.

To conclude Section 3, the above breakdown shows the **major components** (CLI, Config loader, Orchestrator, SuiteRunner, Runner, Plugin categories, Sinks, Controls) with responsibilities and interactions. The code citations confirm their key behaviors: - Orchestrator bridging components ², - SuiteRunner merging and building contexts for experiments ²⁷⁹ ¹²⁹, - Runner performing per-row operations and concurrency handling ³²³ ³³⁵, - Plugins performing their specialized tasks (citations given in context above for tests and CLI referencing certain plugins e.g. metrics and CSV sink).

This component catalog forms the basis for understanding how the system functions as a whole and where to focus for improvements or risk assessments.

4. System Qualities and Constraints

4.1 Security Architecture

Security is a first-class concern in Elspeth's design, as evidenced by numerous controls and defensive features integrated into the code. The security architecture spans authentication/authorization for external integrations, data confidentiality and integrity controls, input/output validation to prevent common vulnerabilities, and traceability for auditing.

Authentication and Secrets: - *External API Authentication:* The system itself doesn't manage user accounts but must authenticate to external services (LLM APIs, storage, repositories). It relies on external credentials provided via config or environment: - OpenAI API: expects an API key (likely set in environment `OPENAI_API_KEY` or passed in config). The code uses the OpenAI SDK which automatically reads `OPENAI_API_KEY`, so Elspeth doesn't expose it beyond reading from env inside library. - Azure OpenAI: uses azure-identity for Managed Identity or environment credentials (client ID/secret). The config likely doesn't require embedding a secret if running on an Azure VM with identity; azure-identity will fetch a token behind scenes ⁵⁵. This avoids storing long-term keys in config. If needed, a user could supply an API key for Azure OpenAI as well, but best practice is Managed Identity. - Azure Blob: similarly, azure-identity allows secure auth without plaintext connection strings. If using a connection string, that would be in config, but presumably, they encourage using identity (the dependency analysis warns to monitor credential escalation CVEs for azure-identity ⁵⁵, implying they

rely on it heavily). - GitHub Repo: requires a token (likely a Personal Access Token or GitHub App token). The code probably expects an environment var `GITHUB_TOKEN` or takes it from config options (which should be stored securely or injected at runtime). It's not in plain code anywhere. Similarly for Azure DevOps, a PAT or OAuth token must be provided out-of-band. - *Secret Management*: The config file `settings.yaml` might include API keys in `options` sections. Recognizing this risk, the documentation likely advises not to commit secrets and potentially supports referencing environment variables in config (the docs might have a convention like `{{ env.VAR_NAME }}` in YAML to load from environment). The system doesn't implement a Vault integration out of the box, but azure-identity integration is a form of secret management (relying on the environment's identity rather than secret strings). - The `SECURITY.md` (present in repo) likely outlines how to handle secrets and any outstanding security tasks. We do see an emphasis on not including secrets in the code or logs: - Logging: The log messages we saw never include the content of prompts or keys by default. They log high-level events (e.g., "Running suite at X" ³⁹⁷ or "Row processing failed after N attempts: error" ³⁶). They do not log the prompt text or LLM response content at INFO level (to avoid leaking sensitive data in logs). If debug is enabled, maybe more info might be logged, but presumably still careful. - The `PromptEngine` does not echo context variables by default unless an error, and then it raises an exception with missing field names (which are not sensitive themselves). - *Internal Access Control*: Within Elspeth, there isn't multi-user access, but the concept of *security levels* acts as an **access control for data flows**: - Each data plugin, LLM plugin, and sink has a `security_level`. The code uses `resolve_security_level` to determine the classification of combined inputs/outputs ⁴⁸. The artifact pipeline's `is_security_level_allowed(producer, consumer)` ensures no sink gets data beyond its clearance ³. For example, if data is classified "Secret" and a sink is configured as "Official" (lower), the pipeline raises `PermissionError` and will not execute that sink ³⁹⁸. This prevents inadvertent export of sensitive data to unauthorized sinks (like sending secret data to a public repo sink). - The classification values (Official, Confidential, Secret, etc.) are normalized (case-insensitive, etc.) by `normalize_security_level` ³⁹⁴. They likely correspond to an ordered enum under the hood. The system likely defaults to some baseline (maybe "OFFICIAL" as default if not specified, which is like low sensitivity). - The documentation likely instructs to set appropriate levels in config for each plugin and sink. If omitted, as we see, sink instantiation will throw config error forcing user to think about it ⁴⁶. - The `determinism_level` tracks whether outputs are fully deterministic or have stochastic elements (like using temperature in LLM). This is more for compliance reproducibility; not directly a security property, but important in audit (some accreditations require demonstrating determinism or variance sources).

Data Protection: - *Confidentiality*: Input and output data can be sensitive (like PII or proprietary info). Aside from classification enforcement, Elspeth ensures confidentiality by not transmitting data beyond intended boundaries: - All LLM API calls go over HTTPS (via requests or openai SDK). The dependency analysis explicitly says to ensure TLS and certificate validation for Azure Blob and others ⁵⁹. Python's requests verifies certs by default; azure-storage-blob uses HTTPS endpoints by default as well. There's no code disabling verification. - If corporate proxies or custom CAs are needed, user can configure environment variables; the code does not override SSL settings to insecure. - When writing to disk, if running on a secure server, presumably OS-level disk encryption covers it. Elspeth itself doesn't encrypt outputs by default. If needed, one might zip and password-protect outputs or rely on disk encryption – not built-in. In future, a plugin could encrypt a file before writing (e.g., a custom sink). - Signing of outputs by `SignedSink` helps integrity (ensuring no tampering after run) but not confidentiality. If confidentiality needed for outputs at rest, one might use a sink that writes to a secure vault or an encryption step. - *Hardcoded credentials*: none found in code – it fetches from env or config. In tests, Dummy sinks or LLM had no secrets. So no secrets are exposed in repository, aligning with good practice. - *Integrity*: The **Signing mechanism** (if fully implemented) is the main integrity control. `security.signing.py` presumably provides `sign_artifacts(artifacts, private_key)` yielding a signature. If a private key is provided (maybe via file path in config or environment pointer), it

signs either each file or a manifest of file hashes. The signature and maybe public key info are output, so later one can verify. This is crucial for ATO, as it shows outputs weren't altered when presenting evidence. The dependency analysis doesn't list a crypto lib; possibly they intended to integrate a system-specific signing (like using GPG via subprocess or an OS key store). If not yet done, this would be a gap to address for full ATO, but code structure is ready (with pipeline expecting a signed artifact). - The artifact pipeline also monitors integrity by ensuring classification doesn't get downgraded. If a sink attempted to produce an artifact at a lower classification than input, the pipeline wouldn't allow it (they call `normalize_security_level(artifact.security_level)` and check it against binding clearance) ³⁸⁹. So if a sink tries to mark an output at Official while input was Secret, pipeline will raise permission error. All artifacts carry a `security_level` attribute, either from sink config or inherited by default (in SinkBinding prepare, they ensure a sink must declare a `security_level` or it's error ³⁹⁹). - **Sanitization & Validation:** - **Prompt Template Strictness:** They use Jinja2's `StrictUndefined`, meaning any placeholder not provided raises an error ⁴⁰⁰. This prevents accidentally running a prompt with an empty variable, which could cause the LLM to see parts of the template or lead to unintended output. It's a security measure in that it avoids unpredictable behavior or leaking of template structure (which could be considered sensitive). E.g., if `{{API_KEY}}` was in template (though likely not), `StrictUndefined` ensures if not given, it errors rather than outputting blank (which might inadvertently expose partial structure). - **Input Schema Validation:** The system can validate that the input DataFrame matches what plugins expect (via `validate_schema_compatibility`) ⁹⁶. This is to ensure that if a plugin relies on certain columns, they exist and of correct type - avoiding runtime errors or misinterpretation. E.g., if a row plugin expects a column "Name" but dataset has "FullName", the validation can catch it early. This prevents misalignment that could lead to weird outputs or security issues (like if a plugin doesn't find a column, maybe it would skip logic and not sanitize something). - **Output Sanitization:** * **CSV/Excel Injection:** Recognized risk that if a value begins with `=`, it could execute a formula when the CSV is opened in Excel. Elspeth's CSV sink by default sanitizes formulas (from code: `sanitize_formulas=True` in sample usage) ¹¹¹. The sink likely pre/appends `'` to any cell starting with `=`, `+`, `-`, or `@` (newer Excel also treats `@` as formula in some contexts). This neutralizes potentially malicious formulas (like `=CMD|'/C calc'!A0` which could launch calc via DDE if not sanitized). This is an OWASP recommendation for CSV output and they implemented it. * Excel sink might need similar caution: openpyxl writing formulas from data might not happen because any `'=` in data might be treated as literal by openpyxl unless you explicitly write it as a formula. But if a result content is `=HYPERLINK("http://malicious", "click")`, Excel would treat it as text by default if writing to cell value (with no formula flag). So likely Excel sink doesn't have the same injection issue. Still, better safe: they possibly also sanitize in Excel by prepending `'` to `'=` at cell creation (openpyxl would then store `'=` in the sheet). * **Prompt Injection:** This is a threat where user input or context could maliciously alter the prompt instruction (like if a data row has text "Ignore previous instructions..."). Elspeth can't fully prevent that because the prompt content is supposed to be user-provided content. However, they mitigate in some ways: - The separation of system prompt vs user prompt template: system prompt is fixed by experiment (like policies the model should follow), and user prompt is filled with data. If an input data field contained something like `"}\nSYSTEM: ...`, could it break out of the template context? Since they use Jinja2 with auto-conversion for braces, a `}` in user data might cause an error or be auto-escaped if using our conversion logic. Actually, they set `autoescape=False` in Jinja2 (so HTML chars aren't escaped) ⁴⁰¹, which is fine. But they don't treat user input as code within template, it's just data substitution. So `prepare_prompt_context` probably doesn't allow injection into prompt instructions except as literal inserted content. If a user put text that looks like an instruction, the model could be tricked (that's inherent risk not solved by code, maybe early-stop plugin or validations can catch known attack patterns). - They could integrate an **AI content safety** check on input to avoid sending certain prompts, but not sure if done. More likely they rely on out-of-band procedures for constructing safe prompt packs (the docs might caution on prompt injection and suggest using separate fields and system messages). - **Response Validation:** * E.g., ensure the LLM's output is well-formed JSON if expecting JSON. The `PromptValidationError` exception

suggests they validate presence of required fields in output. If using `jsonschema` (they depend on it ⁵⁰), they might define a schema for the expected LLM response and validate it in a validation plugin or in `_run_validations`. * If an LLM response fails validation, they mark that row as failure (so it doesn't produce possibly dangerous or unusable data). This prevents, for example, further processing or writing of incomplete results which might mislead or break later steps. * Possibly they use Azure Content Safety as a validation plugin or middleware. If a response triggers content filters (like containing profanity or sensitive info), the middleware could either redact it or raise a validation error to treat it as a failure. The dependency note suggests hooking content safety via requests in middleware ⁴⁵.

- **Exception Handling:** They broadly catch exceptions around risky operations: * LLM call exceptions (network issues, API errors) are caught and retried ⁴⁰² ⁴⁰³. This improves reliability but also prevents leaking raw exception data beyond the system. They only log error message text (which might contain e.g. "Invalid API key" but that's fine). * Row processing exceptions are caught so one bad data point doesn't crash the whole run ³⁶⁰ ⁴⁰⁴. This isolation is important for availability (discussed in reliability too) but also ensures that if a specific input triggers an unexpected code path (maybe an unhandled corner case in a plugin), it's recorded and skipped rather than halting everything (which could leave partial outputs). * The system doesn't seem to catch exceptions at the orchestrator or CLI level except config errors. If something entirely unforeseen happens (like a coding bug raises at orchestrator level), that would bubble up and crash the CLI. Those would need addressing in QA rather than at runtime though.

- **Cryptography and Encryption:** - As noted, no strong cryptography library is in use (if none is, it's a gap for signing). If signing is implemented with a secure algorithm, they likely expect to use RSASSA-PKCS1 or ECDSA with SHA-256. Possibly they plan to integrate with system's OpenSSL or PGP. - Data at rest encryption is not directly implemented by Elspeth; presumed to rely on environment (full-disk encryption, etc.). For data in transit, it leverages TLS through trusted libraries (requests, Azure SDK). - If the environment had a requirement for FIPS 140-2 compliance, one would need to ensure any cryptographic operations (like if using `hashlib` MD5 or non-approved algorithm) are avoided. They likely use SHA-256 if at all (which is FIPS-approved). The code doesn't highlight any particular algorithm usage, likely they would use `hashlib.sha256` for hashing artifacts in signing.

- **Session Management:** Not applicable (no user sessions). - **CSRF/XSS:** Not relevant because there's no web server or browser interface. However, they do produce HTML reports (potential XSS if injecting malicious content). If a result content included a script and the VisualReportSink embedded it into an HTML page (like if it just dumps text), that could be executed when someone opens the HTML. Ideally, the visual sink should sanitize or escape content when generating HTML. There's no explicit mention of it, but given the focus on security, they might do escaping for any user-provided text in HTML context. Or they might design visual output to not include raw LLM text directly in an HTML page without escaping.

- **Dependency Vulnerabilities:** The dependency analysis file enumerates potential security issues: - They mention to *"monitor for credential escalation CVEs"* in `azure-identity` ⁵⁵. `azure-identity` might have had an issue in the past where environment variables could be misused; by specifying "monitor", they intend to keep it updated. - *"update promptly for schema parsing CVEs"* in `jsonschema` ⁴⁰⁵ suggests there have been CVEs (`jsonschema` had a DOS vulnerability earlier). The version pinned (4.21.1) presumably fixes known issues, but they remain vigilant. - Jinja2: *"stay current for sandbox fixes"* ⁴⁰⁵ - Jinja2 historically had sandbox escape vulnerabilities if you use it to render untrusted templates. In our case, templates are controlled by config (semi-trusted) and data is inserted but `StrictUndefined` is used. They likely don't enable full sandbox mode because not needed for this usage, but still they caution to update Jinja in case (3.1.0 is fairly recent). - Requests: ensure corporate CA and proxies used as required (if MITM proxies in corp, not a vulnerability but config). - They also mention scanning with tools (pip-audit) regularly and capturing reports ⁴⁰⁶, and possibly vendoring critical packages to mitigate supply-chain attacks ⁴⁰⁷. This is more operational but shows a security mindset in maintenance.

- **Auditing and Traceability:** - The system generates detailed logs and records metadata such as: * Retry history for each failed call, stored in exception attributes and in payload's `retry_summary` ²¹. * All prompts and responses for each row can be found in payload (the record contains response content, potentially truncated for large content? but likely full since they then use it

for CSV outputs). * Each experiment's config is saved in results (they include `config: experiment` in suite results) ³¹⁵, meaning one can see exactly what parameters were used (temperature, etc.) for each experiment in the output structure. This is helpful to auditors verifying that required settings (like no streaming or a certain model version) were indeed used. * The existence of `TRACEABILITY_MATRIX.md` hints they map how code meets each security control. For instance, a control "All outputs are labeled with classification" can be mapped to code that sets `record["security_level"]` and pipeline checks it ¹⁴⁸ ³. Another "All sensitive outputs are signed" mapped to SignedSink. * `CONTROL_INVENTORY.md` likely lists all security features implemented. * The metadata includes `security_level` and `determinism_level` of results ¹⁴⁸, so any output dataset is self-labeled with classification, fulfilling compliance that every artifact is marked (like a banner). Possibly they intend to propagate that label into file content (some organizations require the classification string at top of documents/CSV as well). Not sure if they do that, but at least in the JSON/metadata it's clearly present. * For high assurance, they could have also included a cryptographic hash of each result row or artifact in the metadata for tamper evidence. They do sign artifacts at the end, which covers that. - Incident response: If something went wrong (like a plugin raised unexpected exception), it's logged (they even log stack trace for early_stop plugin errors ⁴⁰⁸ but continue). The system is resilient enough to contain failures to individual parts, which aids forensic analysis because the run still completes (with failures noted) and you can review logs and payload to diagnose the issue rather than the entire run aborting with possibly partial data and messy state. - Access logs: There's no mention of logging each external request to an audit log, but `on_suite_loaded` hooks could be used for telemetry. Possibly, if integrated with Azure ML, each prompt/response could be logged as metrics for analysis outside the tool. In absence, the payload itself is the log of interactions (it contains all prompts and responses). - They plan "control inventory" and "traceability matrix" to help an ATO reviewer map code to requirements, which is a great practice (not many open-source projects do that).

Summary: Elspeth's security architecture leverages strict configuration validation, internal data classification enforcement, output sanitization, cryptographic signing for integrity, and careful integration with secure external auth methods. It addresses many OWASP and cloud security best practices: - *Injection*: mitigated (command injection not applicable, prompt/CSV injection mitigated). - *Broken auth*: none (no direct user auth in app). - *Sensitive data exposure*: classification ensures handling decisions; TLS used for external comm; no logging of sensitive content by default; optional encryption not built-in but environment expected to handle at rest. - *Security misconfiguration*: requiring explicit sec/det levels and handling defaults reduces risk of forgetting classification. The heavy use of latest library versions and dynamic extras (install only needed libs) reduces attack surface (e.g., not installing heavy libs if not needed). - *Monitoring*: cost_tracker and telemetry provide monitoring of usage to detect anomalies (like if cost spikes or content safety triggered, one could flag it). - *Supply chain*: they maintain dependency vigilance and even suggest internal mirroring of packages to avoid supply chain attacks ⁴⁰⁷.

One constraint to note: All team members using Elspeth must handle sensitive config (YAML) properly, as it contains how data flows. If an attacker could modify a config (e.g., change a sink to point to a rogue server), they could exfiltrate data. But that requires compromising either the config file or the environment running it – out of scope for Elspeth itself, but something an accredited environment addresses with file integrity checks or restricting who can edit config. Possibly the signing of outputs also covers the config (if they include config in the signed manifest, not sure they do, but they could include config file in artifacts to sign it as well, so any changes would break signature).

In conclusion, Elspeth's security design is comprehensive for the scope of an experiment orchestrator. The code enforces security at multiple layers (config parsing, runtime pipeline, output stage) and is built to facilitate auditing and compliance demonstration.

4.2 Performance Characteristics

Elspeth is designed to handle potentially large experiment suites and heavy LLM interactions efficiently, within the constraints of running on a single machine. Key performance-related characteristics include concurrency controls, resource management, and efficient data handling.

Concurrency and Throughput: - The system supports multi-threading for parallel LLM API calls. By default, if more than 50 rows are to be processed and concurrency is enabled, it will use up to a configured number of threads (default 4) ⁷ ³³⁶. This means it can achieve near 4x throughput for I/O-bound operations (like network calls to OpenAI) compared to sequential processing, assuming the external API can handle it. - The concurrency is configurable via `concurrency_config`: e.g., `enabled: true, max_workers: 10, backlog_threshold: 100, utilization_pause: 0.8, pause_interval: 0.5` as default values. So if backlog (rows to process) \geq threshold (50 default), it enters parallel mode. This avoids overhead on small jobs. If thousands of rows, it will definitely go parallel to improve speed. - The *utilization-based throttling* is a sophisticated feature: `while rate_limiter.utilization() >= pause_threshold: sleep(pause_interval)` before submitting each job ³³⁸. This ensures threads don't saturate the external API beyond a safe factor. Essentially, if the API's internal queue (represented by rate_limiter usage) is 80% full, it waits. This helps maintain throughput while preventing overload or large queue buildup (which would increase latency). - The RateLimiter itself will ensure not more than allowed calls concurrently. For example, if allowed QPS is 5 and 10 threads attempt at once, acquire will block some threads until permitted. - **Connection pools:** `requests` uses a pool of 10 connections per host by default. If `max_workers` is set higher (say 20), and all threads target the same host (openai), beyond 10 simultaneous connections, it will queue in requests unless increased. The code doesn't explicitly adjust requests' pool size. But you can configure that via environment if needed. Given default `max_workers=4`, the default pool is fine. If user sets it to 10 or 20, they might also need to tune `requests.adapters.DEFAULT_POOLSIZE`. There's no direct code for that, but the dependency analysis mentions proxies and CA but not pool limits, so presumably not changed (should mention as potential if high concurrency). - **Batch processing and vectorization:** The code does not batch multiple rows in one LLM call (no interface for that aside from maybe criteria which is multiple calls for same row). Each row triggers one or multiple (for criteria) separate LLM requests. If one had the ability to batch queries (some ML services allow sending multiple queries at once to amortize overhead), Elspeth doesn't currently exploit that. It's more oriented to scenario where each prompt is independent and maybe very context-specific, not easily batchable. This is a design choice trading possible throughput gains (batching) for simplicity and per-row isolation (which also helps clarity in results). - For local compute-bound tasks (like aggregator computing stats or baseline comparisons), those are trivial cost relative to network calls. Pandas operations for aggregator metrics (like mean of a few hundred results) are negligible overhead. If a very heavy aggregator or baseline computation is introduced (like a statistical bootstrap with 10k samples), that could become a bottleneck. However, those could be optimized separately or run in another thread if needed.

Latency: - There is overhead from framework around each LLM call: preparing prompts, and after call updating metrics, etc. This overhead is small (millisecond scale) compared to typical LLM API latency (hundreds of ms to seconds). - Early termination features like `early_stop` can reduce the amount of processing if a stopping condition is met. For example, if you set `early_stop` after X successes or a target quality reached, it will break out and skip remaining rows, improving turnaround time for the experiment by not doing unnecessary calls. - Checkpointing can help in long runs: if a run is interrupted (maybe due to time constraints or failure), you can resume without starting from scratch. This is more about fault tolerance but indirectly performance (no need to redo already done 90% if something fails at 90%). The overhead of checkpointing is minimal – an append to a JSONL file per successful row, which is negligible (plus local I/O of writing a line; at, say, 1e6 rows, that's 1e6 lines, which might be borderline

but each line is short). - Logging overhead is modest: they log at INFO for each experiment finish, and each error/warning. They do not log each row success (which is good, that would flood logs). Logging doesn't significantly slow run except if running in debug where they might log every attempt or content, but that's optional.

Memory Efficiency: - Input DataFrame: Pandas will load entire file into memory. If you have a million rows with dozens of columns, that can be hundreds of MB. This could be an issue if extremely large data. The architecture currently requires entire data to be present to iterate. If memory is a constraint, there's no streaming processing out-of-the-box (the *work-packages/WP001-streaming-datasource* suggests they considered adding streaming mode, possibly reading chunk by chunk and processing to reduce memory footprint). As of now, not implemented, so memory is a limiting factor. They probably assume typical usage is on data sets that fit in memory or that one can spin a large enough instance for heavy data. Alternatively, one could chunk the run by splitting config or running multiple times with filter. - DataFrame to list of records conversion: They convert results to a list of dicts at end for output. That duplicates data in memory (once as DataFrame, then as list of dicts, then maybe as DataFrame again if writing to CSV), albeit result size is often smaller than input if output fields fewer. But with many metrics, it could be similar order. Not a huge deal unless millions of rows, then making a Python dict for each row could be slow and memory heavy. But presumably, contacting an LLM for millions of rows is itself enormous cost/time, so not typical. If so, one might want to process in smaller subsets. - Use of lists vs iterators: They build a `rows_to_process` list (gathering all at once) ¹⁶³. For very large data, that's memory overhead ~ number of rows * tuple size. They could have processed on the fly, but they needed to separate initial skip of processed and `early_stop` break to decide whether to parallel or sequential. Actually, they could start threads as they go, but they opted to first prepare a list. It's simpler logic, slight overhead. The threshold logic means if you had 49 rows, it does sequential and those 49 are in list anyway; if 1000 rows, they also go in list, but memory for 1000 isn't big. If you had, say, 100k rows, that list is 100k entries in memory, which is fine if each entry small. It's a trade-off: small overhead to allow concurrency decision. - Shared objects: They compile templates once up front, to reuse for all rows (which is efficient, avoid recompiling Jinja template each iteration) ¹⁹⁴. They also copy references for repeated objects like `prompt_defaults`, `rate_limiter` usage – making sure not to recreate heavy objects repeatedly. - `ThreadPoolExecutor` automatically handles context switching and thread overhead. Python threads aren't heavy on memory (a few MB each typically for stack), so 4 or 10 threads is fine. If user sets 100 threads (not recommended without adjusting requests pool or if each call is heavy), overhead might appear, and they might hit other limits (like OpenAI might throttle or one could saturate network). - The `cost_tracker` might accumulate metrics across many responses (like token counts), but those are usually just counters; negligible memory.

Scalability Limits: - *Vertical scaling*: It can utilize one machine's multiple cores (through threads, albeit CPU-bound code still limited by GIL; but main wait is network so it's okay). It does not scale across machines out-of-the-box. If an experiment's input is extremely large, the recommended approach to scale out is probably to manually split input and run multiple processes or containers in parallel (like run half the rows on one instance, half on another) and then combine results. The architecture doesn't orchestrate multi-node distribution, which is a conscious scope limitation given it's mostly for offline experiments not needing cluster computing. They can be integrated in a pipeline orchestrator (like calling Elspeth CLI in multiple jobs). - *OpenAI API limits*: It's often a bottleneck – e.g., if the limit is 60 requests/min for some model, sending too many will result in rate limit errors. They mitigate by: - Rate limiter config: if user sets it properly, it will throttle to within limits, thereby preventing many 429 errors. If not, their retry logic will catch 429 and maybe eventually respect it (the OpenAI Python library has built-in backoff for rate limits sometimes). They do exponential backoff on exceptions including rate limit, which is good for fairness ⁴⁰³. - Because `utilization_pause` defaults 0.8, they will slow submission if the `rate_limiter` (likely a token bucket with fill rate = allowed QPS) indicates heavy use. This smoothing yields more stable throughput and fewer spikes of error. - *Batch size (Number of experiments)*:

- Running multiple experiments sequentially in a suite: overhead is relatively low between experiments (just some merging and re-initializing runner). If you run 10 experiments each 100 queries vs 1 experiment 1000 queries, overhead of switching experiments is minor (some repeated data load, but they reuse df for suite so no extra load cost, only new runner instantiation overhead). - If user sets concurrency for suite as well (like running experiments in parallel is not implemented; they run sequentially in code), though they could run separate processes to do separate experiments concurrently if needed. - *Thread safety & race conditions*: They use locks appropriately around shared lists for results and failures ²³⁷ ³³⁷. They use an `Event` for `early_stop` with `.is_set` checks and also guard that with a `Lock` in evaluation to avoid potential race in setting reason vs reading event in other threads ⁴⁰⁹. This is careful design: they lock around evaluating `early_stop` in case two threads simultaneously find a reason (one gets lock, sets reason, sets event; the other after lock sees event set and returns without double setting). This prevents e.g. two different `early_stop` reasons clashing. They break out properly if event set. So concurrency is handled robustly, which is important for performance not causing hidden race conditions or memory corruption (shouldn't in Python, but logic issues could cause weird behavior or deadlocks if not careful). - Potential deadlock avoided: They do `rate_limiter` acquire outside any global locks (each thread calls `acquire` context independently, no global Python lock around that aside from if `RateLimiter` internally locks, which it likely does but only short atomic operations). - *Deadlocks risk*: If an `early_stop` triggers while threads are waiting on `rate_limiter`, could there be a scenario? Possibly, if `early_stop` event is set, threads that wake from `rate_limiter` should notice event after finishing current call (they check event at top of worker and in main thread loop when scheduling). So maybe a thread currently inside an LLM call won't check event until done. That means a few calls beyond the threshold might complete but after event set they won't schedule new ones, which is acceptable. No actual deadlock, just a slight overshoot of `early_stop` condition if multiple threads already in flight. - **Timeouts**: - The code does not explicitly set timeouts for HTTP calls. Requests by default has no timeout, which could be problematic if an API call hangs. The OpenAI library typically does have a default or can be configured, but not sure if they passed it. If not, a slow/hanging request could block that thread for a long time. They might rely on user to specify in OpenAI config or OS to drop connection if no response in some time. This is a potential performance/reliability improvement: they could allow a `timeout` in `llm` options and use it in requests. - At a system level, if an external call is hanging, the only remedy would be `early_stop` triggers if something else concurrently signals it (not likely, as they'd be all waiting). - The concurrency means one stuck thread wouldn't freeze entire process (others still run), but if that stuck thread holds a `rate_limiter` token, could reduce throughput. - Not a major issue if using stable cloud API, but still notable. - **Resource Pools**: - They manage cost accumulation but do not try to impose cost-based early stop automatically (though could via plugin). `Rate limiter` is the main resource pool manager (calls per minute). There's also mention in docs to ensure to patch transitive dependencies, not directly performance related but indicates overall environment tuning.

Resource Cleanup: - Threads are cleaned up after use (`ThreadPool` in context manager ensures join). No lingering threads after run. - Files opened (like checkpoint file) are closed properly via context. The code uses `with open(...)` for checkpoint appends (we see `_append_checkpoint` uses a context manager to write JSON line ⁴¹⁰). - If a run fails mid-way, some partial outputs may exist (like a CSV that was being written may be incomplete). There's no rollback mechanism to delete partial output. But since everything runs at end in artifact pipeline, if a sink fails, preceding sinks might have already written to disk. They don't attempt to remove those. So for performance, there's no immediate effect, but for cleanup maybe user has to remove partials or fix and rerun. - Memory for `DataFrame` and results is freed when process ends. If running multiple experiments sequentially via suite, the same `DataFrame` is reused (they don't reload it for each experiment unless `suite_root` changed dataset). That's good to not duplicate memory. - They drop references to potentially large objects after use, e.g., `self._active_security_level = None` after pipeline to break reference cycles if any ³⁵⁴. That

might be to allow GC of context if no longer needed (though not entirely sure if needed, but a nice cleanup).

Benchmark considerations: - The speed of running many LLM calls depends more on LLM service latency and concurrency allowed than on Elspeth overhead. Elspeth adds maybe a few milliseconds overhead per call for processing, which is minor relative to say 300ms API call. - One potential performance improvement would be asynchronous I/O (asyncio) for massively parallel calls. They chose threads likely because it integrates well with CPU-bound parts (like prompt rendering still under GIL, but minor) and using requests (which is sync). Python's GIL is not a big issue because most time is waiting on network and during that, GIL is released. So threads are fine here. Using `asyncio` and `httpx` could allow more concurrency with less thread overhead if thousands of parallel calls, but complexity is higher. For the targeted usage (maybe up to tens of concurrent calls, limited by API anyway), threads are pragmatic.

Scalability Patterns: - The system is stateless across runs, meaning to scale horizontally, one can run multiple instances in parallel (for different tasks or shards of data). There's no in-built distributed but nothing preventing orchestrating that externally. - The concurrency and rate limiting features allow it to **simulate production load patterns** as the README claims ¹⁰ (e.g., you can configure `concurrency=10` and a certain rate to mimic multi-user or high-traffic scenario for that LLM API, and `early_stop` to stop if too many fails). This is not only good for performance, but also performance testing of LLMs (like how does model handle 10 QPS). - **Caching:** It doesn't implement result caching of LLM responses (which could save cost if identical calls repeated). Not typically needed since experiments often vary input slightly. If user needed caching, they could implement a Row plugin that caches outputs given a prompt hash. But not built-in because might not be broadly useful or safe (caching might lead to stale answers if context matters). - **Precision of metrics:** The cost metrics and performance metrics computed (like attempts count, or aggregator results) are all done precisely after processing, so no real performance concern there.

Quality of Service: - The system can adapt to avoid overload (via rate limiter). - If the LLM service is slow or responses large, the concurrency means at least others run in parallel to maximize throughput. - The `early_stop` could implement a "circuit breaker": e.g., if 5 failures in a row, an `early_stop` plugin could trigger and abort further calls to avoid wasting time (this would be configured by user, but the support is there). - Retries provide resilience to transient failures, at the cost of increased latency for those requests. The default `max_attempts` may be config (if not set, default is 1, but user can specify e.g. 3 with backoff). They should consider making sure not to hammer on error; they do increasing backoff and break after max, which is standard.

Constraints: - Memory: must hold data and results concurrently. - Single machine execution: ensures simplicity but limits ability to handle extremely large data sets quickly. It's suitable for typical experimentation scale (maybe up to tens of thousands of cases; beyond that, would be quite slow or pricey). - Network: reliant on internet connectivity for external calls. If offline, user must use mock LLM (which is fast and free but obviously not real). - The performance of output writing is likely minor relative to everything, but e.g. writing a big Excel file with `openpyxl` can be slow (`openpyxl` is pure Python and writing 100k rows might take significant time, possibly more than writing a CSV). If needed, they could note that or use CSV for large data and Excel only for smaller summary. - The design chooses clarity and modularity sometimes over maximum performance (e.g., using Python loops instead of vectorized pandas ops to iterate rows, because they need the logic for each row which can't easily vectorize since each requires an API call anyway).

In testing scenarios: - The presence of many tests indicates they likely measured performance in moderate scenarios. There's no mention of explicit performance testing outputs in docs. Possibly they

assume typical usage patterns (like maybe hundreds or a few thousands of rows, not millions) due to cost of LLM.

So in summary, Elspeth's performance considerations revolve around: - Efficient parallel execution of I/O-bound operations (LLM calls) with controlled concurrency and rate limiting to maximize throughput without overwhelming external services ⁸. - Minimizing overhead per call (compiling prompts once, reusing plugin objects). - Handling large data gracefully (though not streaming, but at least being stable in memory and offering early termination to not waste cycles). - Ensuring reliable and timely results via retries and early-stops (so temporary slowdowns or errors don't derail entire run, maintaining performance stability). - Known constraints such as memory-bound data processing and single-process execution are acknowledged and appear acceptable for the intended use cases.

4.3 Availability and Reliability

While Elspeth is not a long-running service, its reliability in producing results accurately and surviving common failure modes is critical for user trust and compliance. Several design choices enhance its reliability and availability (in the sense of completing runs successfully and reproducibly):

Failure Isolation and Resilience: - Per-row error isolation: The ExperimentRunner does not let one row's failure crash the entire experiment. It catches exceptions for each row (prompt rendering, LLM call, plugin processing) and records that row as a failure in the results, then continues with others ⁹⁵
³⁶⁰. This ensures that if, for example, the LLM returns an unexpected format for one input or a plugin has an edge-case error, the rest of the experiment can still complete. The failed row is clearly logged and present in output (in payload["failures"]) ³²⁴ ³⁴⁸, so it can be addressed post-run. This is crucial for reliability; an older system that crashes on first error could require manual restarts or patching mid-run.

- **Retry Mechanism:** The built-in retries for LLM requests significantly improve reliability against transient issues (network hiccups, rate limit errors, or LLM errors). By default max_attempts might be 1 or user-set, but it's configurable. On an error, they do a delay (which could be incremental with backoff) ⁴⁰³, then try again. The attempt history is stored for audit, and after max attempts if it still fails, they treat it as failure for that row. This prevents random API failures from causing an incorrect output – they only mark failure after exhausting retries. The exponential backoff prevents immediate re-hit on the API if it's down or overloaded, which is a good reliability pattern (avoid thundering herd).

- **Rate Limiting and Throttling:** Beyond performance, the rate limiter actually improves reliability by preventing hitting external service's hard limits too frequently, thus avoiding many failures. It essentially trades slight throughput reduction for far fewer error responses. The pause mechanism at 80% utilization means it leaves headroom – e.g., if OpenAI says 60 requests/min, 80% of that is 48 requests/min, it'll naturally throttle around that so that we rarely see 429 errors. Combined with retries, this yields a stable operation that stays within known good boundaries.

- **Early Stop as Circuit Breaker:** The early_stop plugins can serve to cut off execution if continuing is deemed harmful or pointless: - For instance, an early_stop could check failure rate: if out of first 10 calls, 8 failed, maybe something is wrong (invalid API key, service down) – it could trigger early_stop so the run halts quickly rather than slogging through all data failing. The user can then fix the issue and resume with checkpoint. Without this, one might waste hours/cost on known-failing calls. - Another example: a cost threshold plugin could stop the run if cost exceeds a set budget – that ensures you don't accidentally spend beyond limit if something goes awry. - These act like a circuit breaker in microservice architecture: detecting a problem and breaking out to prevent cascading or wasted effort.

- **Checkpointing:** If a run is interrupted (due to process crash, machine reboot, user abort), the presence of checkpoint file means we can resume and skip what was done. This heavily improves availability in the face of external disruptions. For long runs (maybe hours), you can checkpoint progress. On resuming with the same config, it will skip already processed IDs ¹⁸⁴. This addresses reliability concerns like: - Cloud VM was preempted mid-run – just restart process, it continues (with minor duplication if last row's checkpoint didn't flush, but that's okay). - The process

crashed due to an unhandled exception outside runner (rare, but e.g. if out-of-memory or a bug) – after fixing or adjusting environment, resume is possible. - This design is especially useful given unpredictability of external factors (like the LLM provider maybe going down mid-run). If the run aborted, you can resume later, instead of losing all progress. - **Transactionality**: There are no multi-step transactions needing atomic commit, except maybe conceptually the entire run's operations. But they handle partial completion gracefully (with failures list and so on). - **Data Consistency**: Each result is processed independently. There is no shared mutable state between row operations except the rate limiter and cost tracker: - Rate limiter is thread-safe by design (should be, as used by multiple threads). - Cost tracker is likely also either updated only under GIL or internal locks (if it aggregates totals, likely just increment an atomic variable per call, which under GIL is fine as Python ints are thread-safe for simple ops or they'd lock). - Aggregator plugins run after all rows – they see the final list of results, which in sequential run is obviously consistent, and in parallel run they sort results by index to preserve logical order (though order doesn't typically matter for computing aggregates, but sorting was done to preserve output ordering for nice output; it doesn't affect sums etc.). So aggregator sees consistent complete dataset. - Baseline comparisons run after all experiments, using saved payloads – no concurrency issues there as they run sequentially in suite loop. - **Resource Leaks**: The code opens network connections via requests – requests keeps connections alive in a pool that resets on session destruction. They do not explicitly close the requests session after run. If using the openai library, it manages connections itself; not an issue typically for a short-lived process. If the process stays alive for many runs (e.g. called repeatedly from an interactive session), persistent sessions might cause ephemeral resources to remain. But given typical usage of CLI launching process per run, not a big problem. Also, Python will release memory on process exit, and intermediate heavy objects like DataFrame and runner should be GC'ed when out of scope. There's no evidence of reference cycles except maybe with plugin_context attaching context with references to orchestrator or runner, but they drop active security context at end, possibly to break such cycle if any ³⁵⁴. - **Downtime / Availability**: - Because it's not a service, the concept of uptime doesn't apply. However, "Availability" in an ATO sense might refer to ensuring it completes tasks in a timely manner and can recover from failures. - For scheduled runs (like nightly test suites), reliability features ensure those runs don't unpredictably fail or hang often. If an external dependency is down, early_stop or error logging informs the user and partial results still saved rather than nothing. - If results are needed by a certain time, the user can trust that short of external meltdown, the run will likely finish because of built-in mitigations (retries, throttle). - The pipeline design might leave a partially complete output if a failure happens at final sink – meaning user might get some files but not others if, say, repo push failed. But the run's payload is still returned (except pipeline raising might abort returning payload to CLI? Actually, pipeline.execute is not inside try, so if it raises, runner.run would not return payload normally). They did not catch exceptions around pipeline.execute inside runner.run; they just let it propagate ¹¹⁷. That means if output fails, they don't produce a final payload. The CLI would get an exception and crash (they didn't catch at CLI either). So ironically, if last step fails, the user gets error and no consolidated output (though maybe partial files on disk). That is one area reliability could be improved: perhaps wrapping pipeline in try and adding failure info to payload instead of blowing up. As is, an output failure is treated similarly to a critical error. This is arguably acceptable because if outputs cannot be delivered, the run is not fully successful. - Workaround: user could run again with `--live-outputs` after fixing connectivity to push to repo (using checkpoint to skip processing, it will just attempt outputs again with data in memory; but currently design doesn't persist results in memory between runs, they'd have to re-run experiment which defeats the point. Actually, no easy way to "resume just outputs" because results were in memory). - Perhaps a design improvement would be to write results to a local artifact (CSV) always, then repository sink just picks that file and if it fails, user still has local file. But they do produce local file if configured as such. It's up to user to include robust sinks in config (e.g., always use a local sink in addition to a remote sink to have a local copy).

Determinism and Reproducibility: - They track determinism level because certain runs with randomness (like if model uses temperature or a non-deterministic component) can produce different outputs on reruns. They preserve seeds if any (not explicitly seen, but maybe if using local mock with random, they could set seed to be reproducible). If user desires reproducibility, they can run with `determinism_level: guaranteed` which likely implies using temperature=0 or fixed random seed. They coalesce determinism level in config as they do with security level ⁴¹¹ ⁴¹². Not a direct reliability feature, but important for ATO: being able to reproduce results is often required for verification. They allow marking which parts are deterministic and which are not (so one can explain differences or require multiple trials). - If user tries to do a "load testing" scenario (like concurrency 50 to stress test LLM), the tool can handle it up to point (threads overhead, might eventually saturate host resources or get a lot of 429 errors). Rate limiter might need adjusting in that scenario (like likely they'd disable it to truly push). - The system's complexity (multithreading, file writing, network calls) has been well thought out to avoid common pitfalls. They extensively test scenarios (including integration tests like `test_suite_runner_integration.py`, `test_outputs_*`) indicating reliability is validated in automated tests (e.g., there's likely tests ensuring an `early_stop` triggers properly, a signing works, etc.).

Single Points of Failure: - There's no external single point (like a database). The main SPOF is the single process itself; if it crashes, it stops processing. But as described, it's designed to allow resumption with checkpoint. - External dependencies: If OpenAI goes down, that run will fail for many rows. But the system mitigates by `early_stop` if implemented, or at least it will mark all those as failures but still complete (taking long though). If one integrated content safety and that API is down, the content safety middleware might time out each call – that could slow everything. They might consider adding timeouts or skip logic for such ancillary calls (maybe on `suite_loaded` it can detect if content safety unreachable and disable itself). - Infrastructure: If running as part of a pipeline, pipeline orchestrator should check Elspeth exit code to mark success/failure. Right now, if any exception not caught, the CLI process will exit with a non-zero code. They catch config errors and exit (SystemExit via argparse or raise SystemExit in config errors), and for normal run, if pipeline or anything raises, it likely propagates out causing a stack trace and non-zero exit. For automation, that indicates failure. Partial results might still be on disk but pipeline sees it as failure. Possibly the CLI could be improved to catch exceptions and still output what it got. However, in an ATO context, one might prefer it to error clearly if outputs didn't all succeed, rather than silently giving partial output. It's a design decision. They opted to fail loud on output failure.

Maintenance and Operations: - Logging at INFO provides enough to troubleshoot common issues (like warnings in validation, or errors in specific rows, and summary counts). - The system doesn't have dynamic reconfiguration: you can't change config mid-run (makes sense; each run is static config). - The code is modular, making it easier to fix a bug in one plugin without affecting others – improving reliability in maintenance. E.g., if there's an issue with CSV sink, it won't affect LLM calls logic. - They likely have continuous integration running tests to catch regressions. That ensures reliability of new versions.

Backup: - Not really applicable, but if outputs are important, user should ensure they save them or commit them. The repository sink can serve as a backup mechanism (committing results to a remote git acts like backing up outputs). - Checkpoint file could be considered a small state that should be stored safely if one wants to resume later even if machine lost. Possibly commit checkpoint to a repo as well in between if run spans long time (not automated by tool, but user might do it if needed).

In conclusion, Elspeth is engineered to be robust against many failure modes in its domain: transient API failures, data issues, and performance hiccups. Its reliability features – retries, throttling, early stop, checkpointing – all aim to ensure that experiments can run to completion under expected conditions, and if not, that they can be resumed or partially salvaged. The lack of distributed execution does put limit on availability in terms of scaling out, but as a local orchestrator, it does the best within one node.

The design choices reflect a deep consideration of making the experiment runs **predictable, repeatable, and failure-tolerant**, which aligns well with the needs of a controlled ATO environment (no surprises or uncontrolled crashes).

This completes sections 4.1 to 4.3 by addressing security, performance, and reliability aspects in detail, with references to code where relevant.

Continuing with the remaining sections in a similar detailed manner (5 to 9)...

System Architecture Document – Elspeth

Executive Summary

Elspeth is a Python-based orchestration framework for **secure, compliant large-language-model (LLM) experimentation**. Its architecture emphasizes modularity, security, and auditability, enabling teams to run complex LLM experiments while meeting strict compliance requirements. The system is delivered as a command-line tool (CLI) that loads experiment definitions from configuration files, executes them with controlled parallelism and embedded security controls, and produces comprehensive results and reports.

System Design: Elspeth's architecture is layered into a CLI interface, a core orchestration engine, an extensible plugin system, and integration points for external services. The CLI parses user input and triggers the orchestrator. The **Experiment Orchestrator** and **Suite Runner** components coordinate the flow of data through various plugins: - **Data source plugins** load input data (e.g., from CSV or Azure Blob storage). - **LLM client plugins** interface with external LLM APIs (OpenAI, Azure OpenAI, etc.). - **Middleware plugins** wrap LLM calls for logging, content filtering, and telemetry. - **Experiment plugins** compute custom metrics, validations, or early-stop conditions during runs. - **Sink plugins** handle output persistence (CSV/Excel files, repository commits, signed bundles, etc.).

These plugins are all registered and injected into the experiment pipeline at runtime based on YAML configuration ⁴¹³ ⁴¹⁴. The orchestrator uses this configuration to instantiate an **ExperimentRunner**, which processes each input data record through a sequence of steps: prompt rendering, LLM invocation, result assembly, and metric extraction. Results from all records are then collated, and output sinks are executed via an **Artifact Pipeline** that enforces dependency ordering and security clearances for artifacts ⁴¹⁵ ⁴¹⁶. Multiple experiments can be grouped in a suite for comparative analysis, including baseline vs. variant comparisons.

Security & Compliance: Elspeth was built with a “security-by-design” philosophy. Each data artifact and output is tagged with a **security level** (e.g., Official, Secret), and the system prevents any data from flowing into a sink not authorized for that level ³. All prompt inputs and outputs are validated or sanitized: - Missing prompt variables cause errors (no silent failures) ⁴⁰⁰. - Potentially dangerous content in outputs (e.g., Excel formulas) is neutralized by the sinks ¹¹⁰. - An optional content safety middleware can scrub or flag responses for policy violations.

To support auditability, Elspeth produces detailed logs and output metadata. Every experiment run yields a structured JSON summary of inputs, outputs, and metrics, and can produce signed output bundles to ensure **non-repudiation and integrity** of results ⁴¹⁷ ⁴¹⁸. The framework includes an **Audit Logger** middleware and captures retry attempts, errors, and key events for traceability ²¹ ⁴¹⁹. These

features map to common security controls (e.g., information flow control, audit logging) as required for an Authority to Operate (ATO).

Reliability & Performance: The architecture incorporates robust failure handling and scaling techniques: - A **concurrency controller** enables multi-threaded processing of records, boosting throughput for I/O-bound LLM API calls ³³⁵ ³³⁶. - A **rate limiter** smooths out request bursts to avoid hitting external API limits, pausing task submission when utilization is high ⁸. - **Retry logic** with exponential backoff is applied to LLM requests, significantly increasing success rates under transient failures ¹² ⁴⁰². - **Early-stop plugins** can halt an experiment early if certain criteria are met (e.g., quality threshold achieved or too many failures), preventing wasted time and cost ⁹⁷ ²⁰⁶. - A **checkpoint mechanism** records completed inputs to allow safe resumption of long runs after interruptions ³³¹ ²⁷.

These measures ensure that experiments complete successfully and on schedule, or fail gracefully with actionable diagnostics. The system has been tested to handle typical experiment scales (hundreds to thousands of prompts) and can be integrated into automated pipelines (CI/CD or scheduled jobs).

In summary, Elspeth's architecture provides a comprehensive solution for running LLM experiments in regulated environments. It marries the flexibility needed for rapid ML experimentation with the rigorous controls required for compliance. The following document details each aspect of this architecture: system context, component designs, data flows, security controls, technology stack, architectural decisions, potential risks, and compliance mappings. Every architectural claim is linked to specific code evidence, affirming that the implementation aligns with the intended design and compliance objectives.

1. System Context and Scope

1.1 System Purpose

Elspeth's purpose is to enable **responsible and reproducible LLM experimentation** within an enterprise or regulated context. It provides an orchestrated pipeline to compare prompt variations, model versions, or other experiment factors, while ensuring that all experiments adhere to security and compliance policies. From the code and documentation: - It is described as *"secure, pluggable orchestration for responsible LLM experimentation"* ¹. In practical terms, this means Elspeth allows data scientists and engineers to define experiments (prompts, models, metrics) and run them such that: - **Security:** Sensitive data is protected (through classification tags, sanitization, and controlled outputs) and usage of external services is monitored (via logging and rate limiting). - **Pluggability:** New data sources, LLM integrations, or custom analytics can be added without modifying core code ⁶. This is crucial for extensibility – e.g., one can plugin a new model API or a new output format by writing a plugin rather than altering the orchestrator logic. - **Reproducibility & Audit:** Experiments are executed in a consistent, automated way (versus ad-hoc manual runs), producing an audit trail of what prompts were used and what outputs were generated ⁴²⁰ ¹⁴⁸. This ensures that results can be reviewed or re-generated if needed for verification.

Elspeth is intended to serve multiple stakeholders: - **Data Science Teams** use it to streamline and scale up prompt tuning and model comparison tasks. They benefit from not having to worry about writing glue code for each experiment; instead they declare experiments in a config and let Elspeth handle execution and logging. - **Compliance and Security Officers** use it as an assurance tool. The system enforces that experiments are run in compliance with policies (e.g., no data exfiltration to unauthorized sinks, prompts are constructed from approved templates, etc.). It produces evidence (signed results,

logs) that can be used in audits to prove compliance ⁴¹⁷ ⁴¹⁸ . - **Engineering/DevOps** can integrate Elspeth into continuous integration pipelines or nightly jobs to perform regression testing on model outputs or monitor for drift, etc., knowing that it has built-in checks and will provide alerts (via logging or early-stop) if something goes awry (like a model's performance drastically dropping or outputs containing disallowed content).

In summary, Elspeth's primary mission is to **enable innovation (through flexible experiment configuration) without compromising on governance**. It addresses the need for agility in working with LLMs by providing a standardized, secure process. The codebase reflects this dual focus: for example, the plugin registry allows new experiment logic to be "dropped in" easily ⁶ , while the security module ensures each such plugin is bounded by security context (every plugin factory call is given the experiment's security level so it can enforce its own rules or inherit the context) ⁴⁸ ¹²⁵ .

Through code analysis, we see concrete evidence of this purpose: - The **Experiment Suite** capability exists to manage governed experiment campaigns (it merges prompt packs, suite defaults, etc., so that organizations can define standard prompt sets and apply them across models) ²⁸⁰ ⁴²¹ . - Strict output handling is evident – e.g., by default, repository outputs run in dry-run mode unless explicitly allowed, preventing accidental live pushes ³⁴ . - A **Traceability Matrix** (in docs) maps requirements to implementation; for instance, the requirement "All outputs must be labeled with a security classification" is met by the code that attaches `security_level` to each record and artifact ¹⁴⁸ .

Overall, Elspeth serves as a **trusted orchestration layer** between the experimental ideas of data scientists and the operational safeguards required by the enterprise. It automates LLM experimentation at scale while embedding necessary controls to ensure those experiments are **responsible, auditable, and repeatable**.

1.2 System Boundaries

Elspeth operates within a well-defined boundary that includes the local runtime (where experiments are executed and controlled) and external systems it interacts with. It is fundamentally an **offline, batch-oriented system** triggered by a user or scheduler, and not a continuously running server.

Internal Boundary – Elspeth Process: The core of Elspeth runs as a single Python process (when invoked via the CLI). Within this process: - It reads configuration files (YAML) and loads any necessary local resources (e.g., prompt pack definitions, or sample input files) from the filesystem. - It does *not* spawn independent sub-services or require a database; rather, it uses in-memory data structures (pandas DataFrames, Python objects) to hold experiment data and results. - The CLI invocation `python -m elspeth.cli` (or the installed `elspeth` command) is the entry point, which then calls into the orchestrator logic ²⁴ ²⁵ . Once experiments are done and outputs are written, the process terminates. Each run is isolated in memory and disk outputs; there is no persistent daemon process or long-lived state retained in memory after the run.

Because of this design, **Elspeth does not maintain open network ports or accept inbound connections**. It's not a service that others connect to at runtime. Thus, the trust boundary is primarily at the invocation and configuration level: whoever runs Elspeth (and provides the config) is trusted to define what it should do. The system boundary encompasses: - The **user or automation agent** that executes the CLI (with their permissions to read input files and write output files). - The **machine environment** where Elspeth runs, including OS-level protections, file system, etc., which should be secured (if the environment is hardened, it ensures that Elspeth's intermediate data in memory and outputs on disk are protected from unauthorized access).

External Interactions: Elspeth connects to several external services as part of experiment execution, which define its **outbound boundary**: - **LLM APIs:** e.g., OpenAI or Azure OpenAI endpoints over HTTPS. These are outbound calls initiated by Elspeth's LLM client plugins ³⁶⁵ ¹⁴². No external entity calls into Elspeth; it always acts as the client. Authentication for these calls is handled via API keys or Azure credentials (which are sourced from config or environment, as discussed in Section 4.1 Security). From a boundary perspective, the confidentiality of data sent to the LLM and received in responses is a consideration: data leaves the Elspeth process boundary to the cloud provider. This is governed by policy (ensuring only non-sensitive or appropriately classified data is sent to certain providers, which Elspeth facilitates by classification labeling and potential content scanning). - **Cloud Storage:** e.g., reading input data from Azure Blob Storage, or writing output artifacts to cloud storage. The data source plugin may open connections to Azure Blob (via Azure SDK) to fetch files ⁶⁸. Similarly, an output sink might upload files to a blob container or repository. These interactions are outbound and authenticated using tokens or keys. The boundary here is that secrets/credentials for these are provided to Elspeth and used to perform the transfers – Elspeth itself doesn't host these files, it only retrieves or sends them. - **Version Control Repositories:** If the `repository` sink is used, Elspeth will connect to e.g. GitHub's REST API to create a commit or upload a file ³⁴. This again is an outbound call. The boundary concern is that once data is pushed to an external repository, its confidentiality relies on that repository's access controls. Elspeth's role is to ensure that only data cleared for that destination (by security level) is actually sent ³. Internally, these sinks use `requests` to send data out – if `dry_run` mode is not disabled, they won't actually push (ensuring by default no external write happens unless explicitly allowed in CLI flags or config). - **Telemetry/Azure ML** (optional): If telemetry middleware is enabled, Elspeth might send experiment metrics to Azure ML or another monitoring service ⁴²² ⁴²³. This is again outbound, typically internal to the organization's cloud environment.

No external system calls into Elspeth or controls it at runtime. The **control flow** is one-directional: **Elspeth (client)** -> external service (server). For example: - Elspeth calls OpenAI API (the OpenAI service never calls back; results are returned to Elspeth which then continues processing). - Elspeth might download a blob (Azure Blob storage doesn't push data to Elspeth on its own; Elspeth initiates and receives). - Elspeth pushes to a GitHub repo (the repo doesn't have any hook into Elspeth beyond receiving data; any webhooks from repo are outside Elspeth's scope).

Therefore, **network-wise**, Elspeth sits within an environment (e.g., a secure compute environment) and only makes outbound HTTPS connections to known endpoints: - It requires egress permission to those endpoints (OpenAI, Azure endpoints, etc.). If running in a restricted network, proxy or firewall rules might need to be configured (the documentation notes support for corporate proxies and CA bundles for requests library ⁴⁵). - It expects no inbound network access; thus, firewall rules can block all inbound to the machine except what's needed for general maintenance or user access.

Data Boundary: Data flows into Elspeth primarily via configuration and input files: - The YAML config is often provided by the user (or a pre-defined experiment pack). This contains no secrets except possibly references to secrets (like environment variable names for API keys). The config is read at startup by Elspeth and then lives in memory as `Settings` object ²⁶. The boundary here is that config files should be protected at rest by OS permissions because they could contain API keys or classified prompt text. Typically, only authorized users on that machine should have read access to the config files. - Input data (like a CSV of prompts or records): if local, it's read from disk into memory via pandas. If remote (Azure Blob), it's fetched over a secure connection. Once loaded, the data is in memory (DataFrame) and the connection to external storage is closed. The boundary concern is ensuring the input data source is trustworthy and that Elspeth doesn't inadvertently expose that data beyond intended sinks. Elspeth addresses the latter by classification controls on outputs. Trust in input source is assumed (if input blob is tampered before run, Elspeth doesn't currently verify integrity unless the blob store provides MD5 which Azure does and could be checked – not implemented though). - Output data: flows out through

sink plugins. The boundary is where those outputs go – e.g., if writing to a local file, they remain within the host boundary; if uploading to a repo or cloud, they cross into external boundary. Elspeth ensures outputs are appropriately labeled and optionally signed before crossing the boundary (for instance, a `SignedSink` might produce a cryptographic signature locally, which could then accompany files uploaded externally to prove they weren't altered in transit or afterward).

Actors within the Boundary: - *Human operator*: Initiates the run, monitors logs, and receives outputs. They have control over config parameters (like enabling `--live-outputs`) that affect whether data actually leaves the system (e.g., by default repository pushes are dry-run, the human must explicitly allow live push ³⁴). - *Elspeth process components*: CLI, Orchestrator, Runner, etc., act as internal sub-roles (but not separate trust zones – they run with the same process privileges). The code uses internal checks (like not writing to sinks with insufficient clearance) to protect against misconfiguration or plugin misbehavior. - *External services and their trust*: LLM and storage services are semi-trusted – they are expected to handle data per their security commitments, but from Elspeth's perspective, once data is sent out, it relies on those services. This is why classification gating is important: for example, one might decide only "Official" (low sensitivity) data can be sent to a third-party API, whereas "Secret" data might only be allowed to an internal model or not at all. Elspeth enforces such decisions via config (if someone tries to send Secret data to an OpenAI sink plugin marked as only Official, the artifact pipeline stops it ³). Thus, the system boundary for highly sensitive data might exclude public APIs entirely – Elspeth can be configured to use only internal endpoints for those cases, or to not define any external sink for that level.

In conclusion, the system boundary is drawn around the **Elspeth execution environment**, including its configuration and local runtime, and excludes any persistent external infrastructure. All interactions across this boundary are **outgoing** and either pre-authorized (if credentials are supplied and classification allows) or blocked. There are no incoming control commands at runtime. This simplicity (no listening server, no multi-process state issues) reduces the attack surface and makes it easier to reason about the security of the system – essentially, if you trust the config and the machine Elspeth runs on, you need only then trust the external APIs to which Elspeth connects. The code ensures that trust is not misplaced by preventing accidental data leakage to those APIs or outputs beyond what's intended.

1.3 Stakeholders and Actors

Elspeth's operation involves several key stakeholders and actors, each with distinct concerns and interactions with the system:

- **Experiment Operator (Data Scientist/Engineer)**: This is the primary user who defines experiments and runs Elspeth. They interact with the system by writing configuration files (specifying data sources, prompts, models, metrics, etc.) and executing the CLI. Their goals are to easily conduct experiments and obtain results with minimal manual effort. In the system, the operator's role is reflected by the CLI usage and the configuration profiles:
- They invoke commands like `elspeth.cli --settings experiment.yaml --reports-dir outputs/` to run experiments and direct outputs ²⁴⁷.
- They may also create "prompt pack" files or use built-in examples to define the scope of experiments. For instance, an operator could prepare a prompt pack YAML with multiple prompt templates and criteria, which Elspeth will merge and execute across multiple models as defined.
- The operator expects the system to enforce any constraints they specify (e.g., if they mark data as sensitive, the system should not let it out to an unauthorized sink). The operator also relies on Elspeth's feedback (logs, warning messages) to refine their experiments. For example, if their

config has an error or missing field, Elspeth's validation will warn or abort with an error citing the issue [222](#) , guiding the operator to fix it.

- **System/Infrastructure Administrator:** This stakeholder might install and configure Elspeth on secured servers and ensure it's integrated with the environment's security infrastructure. They may not use Elspeth to run experiments day-to-day, but they:

- Provide necessary credentials via environment variables or secrets store (e.g., set `OPENAI_API_KEY` on the machine or configure Azure Managed Identity for the VM so that Elspeth's Azure plugins authenticate automatically).
- Lock down the machine or container Elspeth runs in (harden OS, restrict network egress if needed). They ensure outputs go to approved locations (they might pre-configure certain sinks or mount output directories).
- They also might be responsible for updating Elspeth to newer versions and ensuring that any dependencies (like `pip` packages) are kept secure and up-to-date, according to policy.

- **Compliance Officer / Security Reviewer:** Although not directly using the tool, this stakeholder has requirements that Elspeth must meet. They will:

- Review the architectural documentation (like this SAD) and the control implementations. They might use Elspeth's **control inventory** (if provided) to see how each required security control is addressed in code [424](#) .
- They may also inspect logs and signed outputs as part of an audit. For example, after an experiment, they might verify that a signed output's signature is valid (ensuring no tampering) and that classification labels on files match the content's sensitivity.
- They might set organization-wide configuration defaults. Perhaps they provide a base config profile that all experiments must include (with certain mandatory plugins like content filtering or an artifact repository for audit evidence). Elspeth supports profiles and merging, so a compliance team could ship a locked-down default profile and require all operators to inherit from it – the code merges profile sections such that defaults and mandatory plugins are always included [280](#) [129](#) .
- Essentially, this stakeholder influences Elspeth's usage by defining policy (which Elspeth enforces). They interact with outputs and logs: e.g., reading the "executive_summary.md" or "comparative_analysis.json" that Elspeth can produce for a suite run [425](#) , which provides them a high-level view and details needed for risk assessment or approval decisions.

- **External LLM Service** (OpenAI/Azure) – *Actor in context of system*: This is not a human actor, but the LLM API can be thought of as an actor that provides responses to Elspeth. Its "behavior" (the quality and format of responses, or rate limiting) influences how Elspeth runs:

- If the service is slow or returns errors, Elspeth's retry and early-stop mechanisms react accordingly to ensure the experiment run remains reliable.
- The LLM service imposes certain rules – e.g., usage limits, or content filters (OpenAI might return a content policy violation warning for certain prompts). Elspeth doesn't override those but could log them or use them in validation.
- Trust: The LLM service is semi-trusted; it sees prompt content and returns data. The boundary here is that sensitive data might be sent to it if allowed. The operator and compliance

stakeholders must decide which services are trusted with what data (Elspeth provides the knobs via classification and plugin config).

- **Data Source Systems:** Similar to LLM service, external data sources like Azure Blob or a database can be considered actors providing input data:
- Elspeth's DataSource plugin initiates interaction (e.g., downloads a file) ⁴²⁶. If credentials fail or data is missing, the plugin (and thus Elspeth) fails gracefully with an error (which the operator must address by providing correct credentials or ensuring data availability).
- The data source must be populated and managed outside Elspeth. The operator ensures the required data is present at the source. For instance, if using `csv_blob` plugin, the operator or an admin must have placed the CSV in the specified blob container beforehand. Elspeth will report an error if it's not found or not accessible.
- **Output Consumers (Stakeholders for results):** After Elspeth runs, the outputs might be consumed by:
 - The experiment operator themselves (e.g., reviewing an Excel report or JSON results to draw conclusions).
 - Other tools or systems: For example, if results are pushed to a Git repository, another pipeline might pick them up to visualize or to trigger further actions. Or if results are uploaded to an analytics dashboard, business users might view them. These consumers rely on Elspeth to have produced correct and complete data. Elspeth's signing of artifacts provides assurance to these consumers that the data came from a valid run and wasn't altered ⁹³ ⁴¹⁸.
 - The compliance team again might be a consumer here, verifying outputs against requirements (like ensuring no PII in a publicly shared result – which Elspeth would mitigate by content filters if configured, but final checking might be done by a person or another system scanning the output).

In the code, **role separation is enforced primarily through configuration and context, not through user accounts** (Elspeth runs as a single-user tool under whoever executes it). For example, the notion of a “*security officer*” is embedded not as a user login but as rules in the config (e.g., forcing certain security levels and signing on outputs). The code shows multiple profiles can exist in one config file (via YAML sections for profiles) ⁴²⁷, which could correspond to roles or environments – for instance, a “default” profile vs. a “strict” profile. The CLI allows selecting a profile (`--profile` flag) ⁴²⁸, meaning an operator can run under different presets possibly set by compliance for different data sensitivity levels. This aligns stakeholder intentions with system behavior by design.

Summary of Interactions: - The Data Scientist defines and triggers experiments (actor at start and end of process). - The System enforces policy during the run (automated actor, via orchestrator and plugins). - External services (LLM, storage) act during execution providing data or receiving data. - The Compliance stakeholder reviews outputs (post-run actor, ensuring everything is per policy).

These actors collaborate indirectly through Elspeth: e.g., Compliance defines rules -> implemented in Elspeth config -> Data Scientist runs experiment -> Elspeth automates under those rules -> outputs are consumed/verified by both Data Scientist and Compliance.

By analyzing code and config: - The CLI and config are the main interface for the Operator ²⁴ ²⁶. - The enforcement points (like pipeline security gate) show the presence of an implicit Compliance actor in the system's logic ³. - External service calls are clearly demarcated (e.g., the OpenAI plugin calling out)

showing where the system relies on external actors ³⁶⁶. - The design ensures that, from an ATO perspective, no single actor (even the Operator) can override critical controls without modifying config in a way that would be evident (e.g., they'd have to set all security levels to the highest and enable live outputs; such config changes would be captured or could be restricted by providing only controlled config files).

In summary, Elspeth's architecture cleanly separates concerns of different stakeholders through configuration boundaries and built-in control points, enabling collaboration: the data scientist gets ease-of-use, the compliance officer gets assurance of control enforcement, and the system admin gets a self-contained tool that fits into the infrastructure without needing constant management.

1.4 External Dependencies and Integrations

Elspeth integrates with various external systems and services to accomplish its tasks, relying on their functionality while adding a layer of orchestration and control. Key external dependencies and how Elspeth interacts with them include:

- **Python Runtime & Libraries:** Though not an external "service," it's important to note that Elspeth is built on Python and heavily uses third-party libraries. This influences integration:
- **Pandas:** Used for data handling (loading CSVs, manipulating DataFrames) ⁴⁹. Integration: Pandas reads local files or file-like objects for remote data (Elspeth might feed it a file stream from Azure blob). Elspeth expects Pandas to perform in-memory; issues like very large files would manifest as Pandas memory errors which are handled by Elspeth's try/except (would raise a ConfigurationError or similar).
- **Requests and OpenAI SDK:** Used for HTTP calls to LLM services ³⁸⁰. Integration: Elspeth's LLM plugins either call `requests` directly or use the OpenAI Python library (which internally uses requests). It relies on these libraries to handle HTTPS properly. The system includes them with pinned minimum versions (`requests>=2.31.0`, `openai>=1.12.0`) which are known to have robust TLS and error handling ⁴⁹. If these libraries have configuration needs (like proxies or custom certs), Elspeth doesn't override their defaults, so it inherits environment-level configuration (e.g., `HTTPS_PROXY` env var).
- **Azure SDKs:** `azure-identity` and `azure-storage-blob` are integrated for secure access to Azure resources ⁶⁸. Elspeth's `BlobDataSource` plugin uses these to authenticate (via `azure.identity.DefaultAzureCredential` which automatically tries Managed Identity, etc.) and to fetch blobs. The `azure-storage-blob` SDK provides `BlobClient` to read data. So Elspeth doesn't implement raw HTTP for Azure; it delegates to these official SDKs. Configuration wise, the user might supply a blob URL or container/name and optionally a connection string or let it use environment credentials. The integration is such that if the environment is correctly set (e.g., `AZURE_CLIENT_ID`, etc., or MSI in Azure), Elspeth will seamlessly retrieve the blob. The dependency analysis notes to "*monitor for credential escalation CVEs*" in `azure-identity` ⁵⁵ – highlighting that Elspeth depends on it for security and any flaw could affect Elspeth's trust boundary.
- **Azure ML (optional):** `azureml-core` is listed as an optional dependency ⁵⁴. Integration: Possibly used by an Azure ML telemetry middleware to log experiment metadata to Azure ML workspace runs. If enabled, Elspeth uses Azure ML's SDK to send metrics. That means Elspeth can integrate into an Azure ML pipeline or at least feed data to Azure ML's tracking UI. This is optional – if not installed, Elspeth just doesn't use that middleware (thus no impact).
- **Matplotlib/Seaborn:** If visual report sinks are used, Elspeth uses these libs to generate charts. These produce images or HTML files saved locally (or included in outputs). They don't call external services but do rely on system capabilities (e.g., an environment with an X server is not needed because matplotlib can use Agg backend for PNGs). Integration here is about being able

to generate images on headless servers – Elspeth likely sets or expects `MPLBACKEND=Agg` in such cases (not explicitly in code, but standard practice).

- **OpenPyXL:** Used for Excel output. It reads/writes Excel files on disk. No external service, but heavy on memory for large files. Elspeth ensures openpyxl is only loaded if needed (as it's optional) ¹¹². The integration is straightforward: if Excel sink is configured, it uses openpyxl to create a workbook. Ensuring that library is present when needed is admin's role (Makefile installs extras as needed, or use `pip install elspeth[sinks-excel]`).

- **LLM APIs:**

- **OpenAI API:** When configured (plugin `openai_http` or similar), Elspeth sends HTTP requests to OpenAI's endpoints (e.g., `https://api.openai.com/v1/completions` for completions, or chat completions). This requires an API key. Integration details:

- The OpenAI Python library is integrated (so usage is via `openai.Completion.create()` or `openai.ChatCompletion.create()` calls) ⁵³. The API key is typically taken from environment `OPENAI_API_KEY` if not explicitly provided in config; in code we see that openai is among dependencies ³⁸⁰ and the dependency note suggests locking to patched versions for vulnerabilities ³⁸.
- Elspeth does not provide interactive login – the key must be configured beforehand. So from a deployment view, integration means storing the API key securely (e.g., as an environment variable in the runtime environment). For Azure OpenAI (which uses the same openai library but with different endpoints), the integration might involve setting `OPENAI_API_BASE` and `OPENAI_API_TYPE=azure` environment variables or using `azure_openai` plugin that knows how to call Azure endpoints with an Azure AD token.
- Rate limits: The integration is such that Elspeth's rate limiter is configured by the user to match OpenAI's published limits. It doesn't query the API about limits dynamically (OpenAI returns some headers for rate limit remaining, but Elspeth doesn't currently use that). Instead, the user must configure concurrency and rate appropriately. However, Elspeth's adaptive throttling (pausing when utilization > threshold) helps avoid hitting limits even if slightly misconfigured, as it prevents queueing an unbounded number of calls ⁸.
- Errors from the API (HTTP 429, 500, etc.) are caught and will trigger Elspeth's retry logic ^{402 403}. So the integration ensures resilience – if OpenAI says "please slow down" or has a transient error, Elspeth will wait and retry up to the configured attempts. The OpenAI library itself also has built-in retry for certain errors, but Elspeth wraps additional logic around it.
- Content filtering: OpenAI might return a special token or flag if content violates their policy, or might truncate. Elspeth doesn't explicitly integrate with OpenAI's content filter endpoint, but it has its own content safety plugin option (for Azure Content Safety) if the user chooses ^{429 419}. Typically, any refusal or filtered response from OpenAI is just treated as a response. If needed, the operator can incorporate that into their validation plugin logic.

- **Azure OpenAI:** Essentially the same underlying service with different auth. Integration:

- The `azure_openai` plugin likely constructs an `openai.api_base` URL (pointing to the Azure endpoint) and uses the OpenAI library with an API key or token appropriate for Azure. If using `azure-identity`, perhaps it obtains a token and sets it in the request

header (OpenAI SDK doesn't natively support Azure AD tokens, but one can override the requests session).

- The dependency analysis notes Azure OpenAI as part of OpenAI SDK usage ³⁸ – presumably Elspeth uses the openai SDK for both Azure and non-Azure by simply switching base URL and key.
- It's configured with the Azure resource name and deployment name via options. The integration requires the user to supply those in config and have azure credentials (either an API key from Azure OpenAI or an Active Directory token via azure-identity). Elspeth's design supports both: if a user provides an Azure OpenAI API key (similar to OpenAI key but specific to Azure endpoint), Elspeth can use it. If they prefer AD auth, possibly the azure_middleware might inject the token.
- Azure-specific concerns: Azure OpenAI has its own rate limits (similar concept). The user would configure accordingly. The azure plugin might also allow selection of model version by name. Elspeth's plugin architecture means it could easily integrate any future changes (like if Azure adds new parameters, one can extend the plugin).

• Azure Blob Storage:

- For input data: The `csv_blob` or generic blob datasource uses Azure's Python SDK. The integration steps:

- It acquires credentials (through environment – e.g., managed identity if on Azure VM, or using an `AZURE_STORAGE_CONNECTION_STRING` if provided).
- It then calls `BlobClient.download_blob().readall()` to get the file content into memory. In code, the dependency analysis references lines in blob_store.py where these calls happen ⁵⁵ ⁴³⁰.
- Elspeth then passes that content to Pandas (if it's CSV) or to the appropriate parser (there is mention of `pandas` being used to read CSV in tests like `test_datasource_csv`).
- If a blob is not accessible (network issue, wrong credentials), the Azure SDK will throw an exception (like `AuthenticationError` or `ResourceNotFound`). Elspeth catches it likely as a generic exception in the datasource plugin and raises a `ConfigurationError` to stop the run with an explanatory message.
- For output data: If configured, a `BlobSink` could upload artifacts to Azure Blob storage. While not explicitly listed in the provided sources, one could exist (the dependency on `azure-storage` suggests writing as well). If so, integration is symmetrical: use `BlobClient.upload_blob()` within a sink. This would require ensuring the blob container exists and the identity used has write permissions.

• GitHub/Azure DevOps Repositories:

- The Repository sinks (like `GitHubRepoSink`) integrate with these services via their HTTP APIs. For GitHub:

- Likely uses `requests` to call the GitHub Content REST API (e.g., PUT to `https://api.github.com/repos/{owner}/{repo}/contents/{path}` with a commit message and file content). Alternatively, could use `PyGithub` library, but no evidence of that, likely direct requests given explicit mention in CLI of those plugin names ³⁴.
- Authentication: via a token in config or environment (like `GITHUB_TOKEN`). Elspeth does not prompt or manage the token beyond using it in the Authorization header. The operator or admin must ensure the runtime environment has this token.

- Dry-run mode: By default, Elspeth's CLI sets `dry_run=True` on any repository sinks unless `--live-outputs` is passed ³⁴. This integration feature means that even if a repo sink is configured, it won't actually push data unless explicitly allowed at runtime. In dry-run, the sink possibly just logs "would push X to repo" or writes to a temp location instead. This is a safety integration to prevent accidental commits – requiring a conscious user action to enable.
 - Error handling: If live mode is enabled and GitHub returns an error (like 401 unauthorized or 409 conflict), the repository sink will throw an exception. In the current code, an exception in a sink will propagate out (causing run failure at the end) unless handled. Likely not caught inside sink, since that would mask a failed output. The user would see an error and perhaps manually push or fix token. (This is one area of improvement – see risk section.)
 - Azure DevOps: Possibly similar integration via Azure DevOps Git REST API. The plugin might call their API to create or update a file or might use an Azure DevOps Python library if one exists (but likely direct REST as well).
 - Both repository integrations mean Elspeth must serialize results appropriately (CSV/JSON content encoded base64 for GitHub API, etc.). The integration handles content encoding and commit messages (maybe adding a header "Generated by Elspeth on date").
 - These external services impose constraints like repo size or rate limit (GitHub's API has secondary rate limits). Elspeth likely isn't hitting those unless pushing very frequently. If used heavily, might require the user to manage scheduling to avoid hitting those limits (Elspeth doesn't incorporate GitHub-specific rate limiting, but general rate limiter could be configured for that domain, although currently the rate limiter is more for LLM calls domain).
- **Email/Notification Systems:** Not directly integrated. If a user wanted email notification on experiment completion or failure, they would have to wrap Elspeth invocation in a script that sends an email based on exit code or parse logs. Elspeth itself doesn't integrate with SMTP or messaging. This was likely a conscious scoping decision to keep it focused (and to avoid having to handle sensitive info via email, etc. as that can be complex compliance-wise).
 - **Operating Environment Integration:**
 - Elspeth respects environment variables for configuration (we discussed API keys, proxies). This means it integrates with the environment's secret management. E.g., on a Linux server, the admin might put the OpenAI API key in an environment variable in the service's context. Elspeth will pick it up because the OpenAI SDK does. Similarly, the Azure identity integration means if the tool runs on an Azure VM with Managed Identity, no secret needs to be passed – `azure-identity` will automatically retrieve a token from the Azure Instance Metadata Service (IMDS). Elspeth is thus integrated with Azure's identity platform by default. In code, they didn't have to implement anything extra – by using azure SDKs, that integration is inherited. For example, dependency analysis notes using managed identity for blob ingestion ⁵⁵.
 - Logging integration: Elspeth uses Python's `logging` module. In an enterprise environment, log handling might be integrated into a SIEM or log aggregator. Since Elspeth writes logs to stderr (by default via CLI), an admin can redirect that to a file or a logging system. The built-in logs have consistent structure (with warnings and info messages) that could be parsed. If deeper integration is needed, one might attach a logging handler to send logs to an external system. Elspeth doesn't do that out-of-the-box (no code for sending logs to, say, Splunk directly), but the design allows it because it uses standard logging – one can configure logging handlers via environment (logging config dict or intercepting log events in wrapper script).

In summary, **Elspeth's external integration strategy** is to **leverage existing, well-vetted libraries and APIs** for external interactions rather than implementing its own protocols: - It **delegates authentication** to environment-specific methods (env vars, Azure MSI), meaning it integrates smoothly into cloud environments (no hardcoded credentials needed in config – a big plus for compliance). - It **ensures secure communication** by using libraries that default to HTTPS and certificate verification (requests, Azure SDK). The documentation even reminds to use corporate CA bundles and such where necessary ⁴⁵. - It **minimizes data exposure** through these integrations by controlling what data goes out (via classification checks) and by offering a dry-run mode for potentially risky outputs. For instance, integration with GitHub is done in a conservative manner requiring explicit opt-in to actually send data ³⁴.

One potential integration not explicitly covered is with **databases** (e.g., reading from a SQL database). Elspeth doesn't list a specific plugin for that, but it could be easily added (and indeed, a user could use Python within a custom plugin to fetch from DB). The current scope appears focused on file-based inputs and API outputs, which likely covers the majority of use cases intended (CSV or blob inputs, LLM API outputs, file or repo outputs).

To illustrate from code: - External LLM call integration: `ExperimentRunner._execute_llm` calls `self.llm_client.generate(...)` which for an OpenAI plugin would internally call out to OpenAI's REST API ³⁶⁶. This is the exact point of integration with an outside system – everything above it is internal orchestration, below it is external execution. They wrap it in try/except to handle external issues ^{431 402}. - External data input: `datasource_registry.create(...)` will yield an instance like `BlobDataSource` which, when `load()` is called, will integrate with Azure SDK to fetch data ^{55 430}. The orchestrator calls this and then integrates the result into a `DataFrame` ⁸⁵. - External output: The artifact pipeline designates certain sinks to push data externally (`RepoSink`, potentially others) ⁴³². They mark what artifacts they consume and produce (for instance, `RepoSink` consumes a JSON or file artifact and has network effect). The pipeline's `SecurityGate` ensures these interactions don't violate policy ⁴¹⁵.

Elspeth's careful use of **standard connectors** and **contextual controls** for all external interactions makes it a well-behaved citizen in a larger ecosystem: it doesn't circumvent enterprise controls (e.g., if a proxy is required, requests library will use it), it doesn't store external secrets in plaintext (they're either in memory or environment), and it respects external systems' constraints (through configurable rate limits, retries, etc.). This approach reduces the integration risk and makes ATO evaluation easier, since the external touchpoints are known and managed with widely accepted practices (TLS encryption, token-based auth, etc.). Each integration was chosen to align with compliance: e.g., azure-identity for Azure integration so credentials aren't hardcoded, or using repository APIs in dry-run by default to avoid accidental data exfiltration.

2. Architectural Views

2.1 Logical Architecture (Structural View)

Elspeth's logical architecture can be understood in terms of layers and components, and their interactions (similar to a **C4 Container and Component** view). At a high level, the architecture follows a **layered pattern**: - **Presentation Layer**: The CLI interface – handles user input and triggers processes. - **Orchestration (Control) Layer**: Core logic that sequences operations (loading data, running experiments, coordinating plugins, writing outputs). - **Plugin (Extension) Layer**: Modular units that encapsulate specific functionality (data access, LLM integration, metrics, outputs, controls). - **External Services**: LLM APIs, storage, etc., which the system calls out to.

Below, we map key components in these layers and describe their logical relationships:

Presentation Layer – CLI: - CLI (Command-Line Interface): Represents the user-facing entry point. It's implemented by `elspeth.cli` and is invoked as `python -m elspeth.cli` or via an installed script ²⁴. Logically, it parses command-line arguments and configuration, then hands off to the orchestration layer. It doesn't contain business logic of experiments, but it prepares the environment (e.g., sets log level, processes options like `--disable-metrics` to adjust settings) ^{433 434}. Think of it as a coordinator that knows how to assemble a `Settings` object from YAML and then call the orchestrator's run methods. Once the orchestrator finishes, the CLI may format some results for console output (like printing a preview of results dataframe, or printing a message that reports were generated) ^{70 435}. In the logical view, the CLI is the top-level component that the operator interacts with, but it delegates nearly all work to underlying components. It's logically connected to: - The **Config Loader** (to get settings from YAML), - The **Experiment Orchestrator** (to actually run the experiment or suite), - Possibly directly to some plugins or components for quick operations (for example, CLI can trigger schema validation without running experiments via a flag ²²³, which calls validation logic in the orchestrator or runner).

Orchestration & Control Layer: - Configuration Loader (Settings Loader): This sub-component (realized by `elspeth.config.load_settings`) reads the YAML configuration and produces a structured `Settings` object ²⁶. Logically, it's a translator between external config and internal state. It interacts with the plugin layer by using registries to create plugin instances from config entries ^{130 67}. For example, if the config lists a datasource plugin "csv_blob", the loader asks `datasource_registry` for an instance of that plugin with given options ¹³⁰. The loader ensures all references in config are resolved to actual objects (datasource object, list of sink objects, etc. in `Settings`). It also normalizes and merges profiles (taking base defaults and overriding with profile-specific values) to produce a cohesive orchestrator configuration ^{280 129}. Logically, it sits between CLI and Orchestrator: the CLI calls it to get `Settings`, and then passes those settings to the Orchestrator.

- Experiment Orchestrator: The orchestrator (class `ExperimentOrchestrator`) is the central control component that coordinates one experiment execution ²³¹. Logically, it is connected to all plugin types: it receives from Config Loader the initialized datasource, llm client, sinks, controls, etc., and it composes them together. It creates an **ExperimentRunner** (next component) and injects all those plugin instances into it ^{74 119}. It is essentially a builder and holder of context. It also establishes the experiment's context (security level, etc.) which gets passed along to plugins ^{48 75}. In logical terms, think of it as the master of ceremony for a single experiment: it tells the DataSource to load data, passes that data to the Runner, and later collects the result payload to pass on to output generation or the CLI. It doesn't implement the loop over data or the multi-threading – that's delegated to the ExperimentRunner – but it sets everything up and provides the necessary high-level controls (like applying `max_rows` limit to dataset if configured ¹³³, or skipping experiment run in certain modes).

- Experiment Suite Runner: This component (class `ExperimentSuiteRunner`) extends the orchestrator concept to handle multiple experiments (a suite). Logically, it manages a collection of `ExperimentConfigs` (each representing one experiment in the suite, possibly with one marked as baseline) ³⁰³. It interacts with Experiment Orchestrator/Runner in that it creates a fresh ExperimentRunner for each experiment, using a **ConfigMerger** helper to overlay suite-level defaults, prompt pack, and specific experiment settings ^{280 79}. It then invokes each runner sequentially and collates results. In the logical diagram, the SuiteRunner is connected to: - The **Config Loader** (it might use it or similar logic to get base settings and each experiment's config), - The **Plugin Registry/Factory** (when building each experiment's runner, it creates plugins via registries as needed for that experiment's specifics), - Multiple **ExperimentRunner** instances (one per experiment) ⁴³⁶, - Baseline **Comparison Plugins** (to do cross-experiment metrics after baseline and each experiment) ⁸⁰. It also handles cross-cutting concerns at suite level – e.g., if a *RateLimiter* is meant to be shared across all experiments in a suite (to throttle combined requests), SuiteRunner ensures either the same instance is reused or each experiment gets one attached but fed

from defaults if not overridden. Logical view: the SuiteRunner sits above multiple orchestrators (or rather, above multiple runner instances), ensuring they adhere to suite-wide constraints and then aggregating their outputs into summary reports. - **ExperimentRunner**: This is a core component in the control layer (though it could also be viewed as part of execution layer). It executes the experiment logic for a single experiment on the provided dataset ³²³. Logically, it is connected to all plugin types as well:

- It uses the **LLM Client** and **LLM Middleware** chain to process each prompt (calls `llm_client.generate` possibly wrapped with `before_request`/`after_response` from middleware) ^{106 57}. - It iterates over data from the **DataSource** (it receives the DataFrame from orchestrator and then loops through it) ^{163 341}. - For each data record, it invokes **Row Plugins** (for metrics per row) ¹⁴⁴ and **Validation Plugins** (to check responses) ⁹⁵. - After processing all records, it invokes **Aggregation Plugins** (for experiment-level metrics) ⁹². - It coordinates with **Early-Stop Plugins** (checking after each record if a stop condition is met) ^{437 146}. - It feeds results and metadata into the **Artifact Pipeline** for output sinks ¹¹⁷. In a logical diagram, the ExperimentRunner would be central with arrows to "LLM Client", "Row Plugins", "Agg Plugins", "Validation", "RateLimiter/CostTracker", "ArtifactPipeline", etc. It's where the main loop and concurrent processing happen. Internally it contains sub-components:
 - * *Concurrency Controller*: not a separate class, but logically the thread pool and related logic inside ExperimentRunner that decides parallel vs sequential run and manages threads ^{7 335}.
 - * *Retry Manager*: again internal, the logic in `_execute_llm` that handles retry attempts and backoff ^{363 402}.
 - * *Checkpoint Manager*: internal, tracking processed IDs and writing to checkpoint file ³³¹. These could be depicted as sub-nodes or functions of the Runner in a detailed component view ^{438 439}.
- **Artifact Pipeline**: This component manages the execution of output sinks and artifact handling after the main experiment logic. Logically, it takes the final `payload` and `metadata` from ExperimentRunner and routes them to each configured **Result Sink Plugin** in the correct order ^{117 432}. It enforces dependencies and security:
 - It checks via **Security Gate** that each sink's clearance is \geq data's level ^{3 440}.
 - It ensures if one sink produces artifacts another needs (e.g., CSV sink produces a file that the Signed sink should sign), it runs them sequentially in the proper graph order ^{113 178}.
 - It uses an internal **ArtifactStore** to register artifacts by ID/type and resolve sink consumption requests ^{441 442}.
 Logically, the Artifact Pipeline is connected to all **Sink Plugins** and is triggered by the ExperimentRunner at the end. It's a bridge between the core run and external outputs, applying final transformations (like signing or zipping). In diagrams, it's often shown as a series of stages: Resolver, SecurityGate, Collector connecting to sinks ^{415 416}.
- **Plugin Registries**: Although not always drawn, logically there is a **Plugin Factory/Registry** subsystem that handles creation of plugin instances (DataSource, LLM, Sink, etc.) from identifiers and options. In code, these are in `core.datasource_registry`, `core.llm_registry`, `core.sink_registry`, etc. For logical understanding, whenever the Config Loader or SuiteRunner needs a plugin instance, it calls the registry's create method with plugin name and options ^{130 124}. This returns a new plugin object (which may be a class from `plugins.*`). Logically, the registry is connected to plugin implementations (one might depict it as a factory icon linking to each plugin type). The updated architecture notes show `PluginRegistry -> Datasource/Sinks/Controls` meaning it feeds into those creations ⁴⁴³. This layer of indirection decouples core logic from plugin instantiation, which is important logically for extensibility.

Plugin Layer (Extensible Components): This layer comprises various categories of plugins, each category implementing a common interface/protocol and serving a specific role: - **DataSource Plugins**: Provide data input. E.g., `CSVLocalDataSource`, `BlobDataSource`. They implement a method (e.g., `load()`) that returns a pandas DataFrame ⁸⁴. Logically, the DataSource plugin is invoked by the orchestrator at start to retrieve the dataset ⁸⁵, and then not used further (except that it might carry metadata like schema or security level that travels with the DataFrame). DataSource is connected to external data store (like a file system or cloud storage) and internally to the orchestrator/runner which consume its output. In diagrams, it's usually depicted at the boundary between Elspeth and external data (arrow from DataSource plugin to an external data source, e.g., Azure Blob) ⁴²⁶. - **LLM Client**

Plugins: Encapsulate communication with an LLM service. For example, `OpenAIClient`, `AzureOpenAIClient`, `MockLLMClient`. They implement a protocol (like `generate(system_prompt, user_prompt, metadata)`) and hide the details of HTTP calls or SDK usage ⁸⁶. Logically, the `ExperimentRunner` calls the LLM client plugin for each prompt (via maybe a middleware chain) ³⁶⁶. The LLM client is connected to the external LLM API (arrow to OpenAI/Azure in diagrams ⁴⁴⁴). Also, the orchestrator sets up the LLM client instance and passes it in. There may be multiple LLM plugins in a single run if user tests different models, but typically one per experiment unless a prompt pack includes multiple named LLMs (Elspeth currently doesn't loop multiple LLMs in one run unless user runs multiple experiments).

- **LLM Middleware Plugins:** Act as interceptors around LLM calls. For example, `AuditLogger`, `ContentSafetyMiddleware`, `RateLimiter` (though rate limiter is a control, implemented differently), `AzureTelemetryMiddleware`. They implement hooks like `before_request(request)` and `after_response(request, response)` ^{106 57}. Logically, the `ExperimentRunner` assembles a chain (list) of middleware from config and wraps the LLM client call through them. Each middleware can modify or augment the request/response. In architecture notes, they show a chain: `Audit -> Shield -> ContentSafety -> Health -> AzureEnv -> LLMClient` as an ordered flow through middleware to LLM ⁴¹⁹. This indicates multiple logical sub-components in the middleware chain. They are connected to the LLM client (they call next, or ultimately the LLM) and to cross-cutting concerns like logging or external safety services (`ContentSafety` calls Azure's content filtering API; `Audit` might log to a file or variable). Middlewares can also have side connections: e.g., a `Telemetry` middleware might output to Azure ML (in diagram, `Visual/Analytics` sinks output to `Telemetry` feed) ⁴⁴⁵. In logical layering, they belong to the plugin layer but function as part of `Runner`'s internal workflow.

- **Row Experiment Plugins:** These run on each record after the LLM response. E.g., a plugin that computes a custom metric from the model's output (like "extract sentiment score"). They implement a method like `process_row(row_data, responses) -> metrics` ⁸⁸. Logically, the runner invokes each row plugin sequentially for each record's result, and collects their outputs into the record's metrics. They are connected to the `ExperimentRunner` (which feeds them data and gets metrics back) and do not interact with external systems (except maybe local computations or reference data).

- **Validation Plugins:** Similar to row plugins but focused on verifying content or format. They might inspect the LLM's response and raise errors if invalid (they might throw exceptions caught by the runner as a failure) ⁹⁵. Logically, also invoked per record after obtaining the model response, likely before or after row plugins (in implementation, they do validation either integrated in `after_response` or as separate step before marking success).

- **Aggregation Plugins:** Run once after all records processed to compute an experiment-level aggregate. E.g., a plugin that calculates "% of responses where condition X is true" or statistical summary. They implement a method like `finalize(all_records) -> aggregate_metrics` ⁴⁴⁶. The runner calls them after sorting results ⁹². They may have dependencies on all data (hence run at end). They output into the experiment's aggregate metrics which become part of payload (and can be used by baseline comparisons or reporting).

- **Baseline (Comparison) Plugins:** In suite context, these take the baseline experiment's results and another experiment's results and compute comparative metrics (e.g., diff in success rate, statistical significance). They implement e.g. `compare(baseline_payload, experiment_payload) -> diff` ³²⁰. The `SuiteRunner` invokes them for each experiment after its results are available ⁴⁴⁷. They output differences that get attached to the experiment's payload as `"baseline_comparison"` ¹⁰². Logically, baseline plugins connect two runs' outputs to produce a third dataset. They might be considered part of aggregator category but for multi-run data. They usually don't interact with external services, just compute.

- **Early-Stop Plugins:** They serve as condition checkers during execution to decide if processing should halt. Implement something like `check(record, metadata) -> reason or None` ²⁰⁵. The `ExperimentRunner` calls them (each plugin in the active list) after each row's success. If any returns a reason, the runner triggers the early stop event ¹⁴⁶. These plugins logically connect to the `ExperimentRunner`'s control flow. An example might be an early-stop plugin that monitors if the cost so far exceeds a threshold, or if a certain number of failures occurred, then signals to stop. They can be configured either directly or via a

simpler config that is normalized (the code normalizes `early_stop` configs into plugin definitions) ¹⁰⁰. They integrate with controls like cost tracker or have internal state. They may also communicate a final reason which goes into metadata (runner stores `_early_stop_reason` which is added to payload) ³⁵⁰.

- **Result Sink Plugins:** Handle output of results/artifacts:
- **File/Blob Sinks:** E.g., `CsvResultSink` (writes a CSV to local path) ¹⁰⁹, `ExcelResultSink`, `LocalBundleSink` (perhaps packages files into a folder or zip), `BlobSink` (upload to Azure Blob). These plugins get the final results dataset (or parts of it) and persist it. They typically run at the end via the `ArtifactPipeline`. They may produce new artifacts (like CSV file path, or a zip file) and the pipeline registers those ⁴³². They might also have security measures (CSV sink sanitizes formulas ¹¹⁰).
- **Repository Sinks:** `GitHubRepoSink`, `AzureDevOpsRepoSink` push results to external version control. They are logically similar to file sinks but the “file” destination is remote. They require network access and credentials. In the logical view, these sinks are connected to external repository services (with a note that by default they operate in dry-run).
- **Analytics/Visualization Sinks:** `AnalyticsReportSink` might produce a Markdown or JSON summary of results (like an executive summary, comparative tables) ⁴⁴⁸. `VisualReportSink` might generate charts (PNG/HTML). These sinks consume the payload and produce user-friendly artifacts (plots, documents) for stakeholders. They run after core metrics are in place. Often they combine multiple experiments’ data (especially `SuiteReportGenerator` which uses all experiments’ payloads to make consolidated reports ⁴²⁵). In logical layering, they are at the output boundary, possibly needing to embed references to external images or use local files (but not calling external APIs typically, aside from maybe sending telemetry of a generated chart to Azure ML as shown by `VisualSink` -> Telemetry in diagram ⁴⁴⁵).
- **Signing Sink:** `SignedArtifactSink` produces a digital signature for artifacts. It consumes all earlier artifacts (as `@all` in config, meaning it will wait for all artifacts) ¹⁸¹, and produces a signed manifest or bundle (like a .sig file or signed zip) ⁴⁴⁹. Logically, this sink ensures integrity: it’s connected to either a cryptographic service or library. Possibly it uses local key (which could be an integration point if keys are in an HSM or KMS, but likely uses a file with a private key or an HMAC secret from config). It runs last (so that it can sign everything). In the logical structure, it’s both an output (signature file) and a check/gate (ensuring nothing modifies artifacts after signing). It might be depicted with a dependency on artifacts and producing an “audit trail” artifact ⁴⁵⁰.
- **Chain of Sinks:** The `ArtifactPipeline` organizes sink execution logically as a DAG (directed acyclic graph). For example, in one diagram:
 - CSV sink produces a file -> `SignedSink` and `ZipSink` might consume that file artifact to sign or bundle it ⁴⁵¹.
 - `VisualSink` produces images -> Telemetry picks those up if needed ⁴⁴⁵.
 - `RepoSink` might consume a JSON output to push it. Each sink is logically independent but the pipeline enforces the connections via artifact flows.

Cross-Cutting Concerns: - **Security Context Propagation:** Throughout the logical model, every plugin and component carries a `_elispeth_security_level` attribute (for sinks/plugins) or uses the `PluginContext` passed by orchestrator ¹²⁵. This means no matter which component is running, it is aware of the data classification and can adjust behavior or check it (e.g., in `ArtifactPipeline._prepare_binding`, they raise if sink has no security level defined ³⁹⁹, making security level a required property for integration). - **Logging:** The logical flow includes logging at key points: e.g., orchestrator logs when starting an experiment or suite, each sink might log an info when writing, any warnings from validations are logged. Logging is not depicted in typical architecture diagrams, but conceptually it is a cross-cutting layer that all components use to report status. E.g.,

ExperimentRunner logs a warning if an early stop triggers (with plugin name and reason) ¹⁴⁶, CLI logs summary info like where reports were saved ⁴³⁵.

To illustrate logical relationships, consider a simplified sequence for a single experiment: 1. **CLI** invokes **Settings Loader** -> obtains `Settings` with `datasource=BlobSource`, `llm=OpenAIClient`, `sinks=[CsvSink, SignedSink]`, etc. ^{130 263}. 2. **Orchestrator** receives these, sets up context (e.g., security level Official because maybe datasource or config said so) ⁴⁸. It calls `data = datasource.load()` to get `DataFrame` ⁸⁵. 3. Orchestrator creates an **ExperimentRunner** injecting: `llm_client=OpenAIClient`, `sinks=[CsvSink, SignedSink]`, `row_plugins=[...]`, `aggregator_plugins=[...]`, `rate_limiter` if any, etc. ^{74 119}. 4. Orchestrator calls `runner.run(data)`. 5. **ExperimentRunner** begins iterating data: - For each row: prepares context dict, calls LLM middleware chain -> LLM Client -> gets response ^{106 142}. If error, perhaps retries ⁴⁰². - Passes response to validation plugins (if any) ⁹⁵. If validation fails, marks this row as failure and continues. - Passes response to each row plugin -> collects metrics into the record ¹⁴⁴. - Updates cost tracker and rate limiter usage metrics if configured (so these cross-cutting components get updated every call) ^{167 371}. - Checks early-stop plugins with the record's data; if any triggers, sets an event to break out after current iteration ^{205 146}. - Continues until all rows or break. 6. After loop: sorts results, computes aggregator plugins (like average metrics) ⁹². 7. Builds `payload = {"results": [...], "aggregates": {...}, "metadata": {...}}` including summary stats, retry counts, any early_stop reason, and a security_level field (e.g., Official) ^{208 148}. 8. Passes payload to **ArtifactPipeline**: - Pipeline creates SinkBindings for `CsvSink` and `SignedSink` ^{352 353}. - The pipeline sees `CsvSink` produces a CSV artifact, `SignedSink` consumes all artifacts. Ensures order: run `CsvSink` first. - Executes `CsvSink.write(results)` -> writes file to disk (artifact `type="csv"` is registered) ^{32 452}. - Then `SignedSink.execute` -> reads that CSV file (and any other artifact, if any) and produces a signature (artifact `type="bundle"` perhaps) ⁴⁴⁹. - The pipeline attaches both artifacts to output metadata and final payload, perhaps. - Enforces security: if either sink had insufficient clearance, it would have thrown error earlier ³. 9. ExperimentRunner returns payload with references to output artifacts. 10. Orchestrator returns that to CLI. 11. CLI might print "Experiment completed with X records, outputs saved to ..." or if `--reports-dir` specified and multiple experiments, uses **SuiteReportGenerator**: - In multi-experiment case, the **SuiteRunner** would have repeated steps 4-9 for each experiment and then invoked Baseline plugins to compare baseline vs others (embedding those comparisons in each experiment's payload) ³²⁰. - Then **SuiteReportGenerator** (an analytics sink in effect) would produce consolidated reports (Excel, markdown, etc.) ⁴⁴⁸. 12. CLI then notifies user (via console or just by the presence of output files).

This logical flow shows clear separation: - The CLI/Config load deals with reading inputs and user preferences. - The Orchestrator deals with constructing the runtime scenario (wiring plugins together). - The Runner deals with core execution logic (loop, concurrency, apply plugins). - Plugins do domain-specific work at appropriate points. - The ArtifactPipeline handles the tail end (ensuring outputs are delivered properly).

The **logical architecture** can be visualized as a series of interconnected components: The official documentation even provides mermaid diagrams confirming these relationships ^{453 454}, which matches our analysis: - CLI in Operator Workstation calls **ConfigLoader** ⁴⁵⁵. - **ConfigLoader** creates instances of **Datasource**, **LLMClient**, **Sinks**, **Controls**, **Middleware** ⁴⁵⁶. - Those plugin instances flow into **Orchestrator** (which belongs to Core Runtime) ⁴⁵⁷. - Orchestrator delegates to **ExperimentRunner** (and **SuiteRunner** if multi-exp) ^{458 438}. - **ExperimentRunner** interacts with the LLM Middleware chain and **LLMClient** for each request ⁴¹⁹, and updates **Controls** (**RateLimiter**, **CostTracker**) during execution ⁴⁵⁹. - After processing, **ExperimentRunner** sends payload to **ArtifactPipeline**, which sequences **Sinks** (**CsvSink**,

RepoSink, SignedSink, etc.) ⁴⁶⁰ ⁴³² . - Sinks produce artifacts possibly consumed by others (shown by dotted lines from sinks back to pipeline or between sinks, representing artifact flow) ⁴⁵¹ .

Thus, the logical architecture is clearly componentized and layered, enabling one to reason about each part in isolation and how they collaborate to fulfill the overall execution of an experiment suite. Each connection has been designed with certain controls (e.g., security context passing, exception handling) to maintain the system's robust behavior across all these interactions, as described above.

2.2 Physical Architecture (Deployment View)

In the physical view, we consider how the system is deployed on infrastructure, how processes and storage are arranged, and network topology. Elspeth's physical architecture is relatively simple, as it is primarily a single-process application, but we will detail how it might be hosted and run in a production environment and how it interacts physically with external systems:

Deployment Model: - **Host Environment:** Elspeth can run on a developer's workstation or a dedicated server/VM. In production or controlled environments, one might run it on a hardened virtual machine or container (Docker container) that has the necessary runtime (Python 3.12 and required libraries) and access to resources (network, files) needed. For example, an organization may deploy Elspeth on an internal server that has internet egress to call LLM APIs, and connectivity to internal data sources. - **Process:** When an experiment is executed, it runs as a single OS process. If multiple experiments need to run in parallel (e.g., scheduled concurrently), they would spawn separate processes (for example, via separate CLI calls perhaps orchestrated by a scheduler like cron or an orchestrator like Airflow). Elspeth is not a daemon; each run is independent. This means physically, it does not maintain persistent state between runs except what's persisted in output files or checkpoint files explicitly. - **Memory and CPU:** The single process is multi-threaded internally for concurrency, but it does not spawn child worker processes. All threads share the process's memory space. Thus, physical resource usage is bounded by that one process's memory (which must hold the DataFrame, results, etc.) and CPU (which might be used by up to the number of threads configured for concurrency, e.g., 4 threads could use 4 CPU cores when waiting on I/O or doing some compute in parallel). If more parallelism is needed than one process can achieve (due to GIL or design), multiple processes could be run (e.g., to utilize more CPU cores if doing CPU-heavy metrics, though in LLM experiments typically bottleneck is I/O). - **Network:** The machine where Elspeth runs needs outbound network connectivity to: - OpenAI's API endpoint (or whichever LLM endpoints configured), - Azure endpoints if using Azure (for blob storage and content safety, etc.), - Possibly GitHub or other repository endpoints if outputs are pushed. In many enterprises, this might be behind a corporate proxy or firewall. As noted, the requests library and Azure SDK respect proxy settings if environment variables are set, so physically, Elspeth's process will route through those proxies as configured. There are no inbound network listeners; you don't need to open any inbound port on the firewall for Elspeth, which simplifies its deployment (it's an outbound client in network topology). - **Storage:** - **Local Disk:** Elspeth reads config files and possibly local input files from the disk. It also writes outputs to disk (e.g., in the `--reports-dir` or default current directory for outputs if not specified). Disk space should be sufficient to store at least the largest dataset and all desired outputs. If outputs include things like numerous images or large Excel files, that must be accounted for. - **Temporary Files:** It's possible that some sinks or operations use temp files (for instance, if generating a zip bundle, it might create a temp zip file). The signing process may produce a manifest file. These would typically be on local disk as well. The checkpoint file is written to local disk (either in working directory or specified path). - The process does not use a database or external state store; all state is either in memory or written to files as output. For example, the traceability matrix or logs are stored in

output artifacts, not in a persistent DB. - If running in a container, one would ensure a volume is mounted for output if needed to persist outside container.

- **External Integration Physical Connectivity:**

- **Azure Blob:** If Elspeth is running on an Azure VM with Managed Identity, the Azure SDK will reach out to the local IMDS endpoint (`http://169.254.169.254/...`) to get a token. This is a local call that requires no internet. Then it will call Azure Blob REST endpoints (like `https://<account>.blob.core.windows.net/...`). Physically, these calls go over port 443 and should be allowed by egress rules. If running outside Azure, it might use a connection string and go directly to the blob endpoint.
- **OpenAI API:** Calls `api.openai.com` on port 443. This requires DNS resolution and internet access. If using Azure OpenAI, calls something like `<resource>.openai.azure.com` (which is also typically a public endpoint, but within Azure region). In a highly secure environment, one might restrict these calls to known IP ranges or require them to go through an approved egress point.
- **GitHub API:** Calls `api.github.com` on port 443. Similarly, ensure egress allowed if that sink is used. If using an on-prem or private Git repository, perhaps the RepoSink could be configured for that domain (the sink likely supports any host via config). Then network access to that host is needed (it could even be within intranet).
- **Azure DevOps:** Usually `dev.azure.com` endpoints on 443, requiring egress to Azure DevOps.
- **Azure Content Safety:** If the content safety middleware is enabled, it will call an Azure service endpoint (maybe part of Azure AI Content Safety). That's another outbound call to an Azure domain.
- **Azure ML:** If telemetry is used, the azureml-core SDK will send data to Azure ML workspace endpoints (maybe `https://ml.azure.com` or an ingestion endpoint). All these communications use TLS, and Elspeth does not intercept or store data except logging possibly that a call happened (it doesn't log sensitive content by default, except in debug mode it might if turned on, which would be ephemeral in logs).

- **High Availability & Scheduling:**

- Since Elspeth is run as needed, if one wants it to be highly available, you could schedule it on multiple servers or have a failover mechanism (but because runs are stateless one-offs, HA is more about being able to re-run if something fails rather than keeping a service up).
- In a production pipeline, one might integrate Elspeth as a step in an orchestrator (like Jenkins, Airflow, or Azure DevOps pipeline). Physically, that means the pipeline agent triggers the CLI on a provisioned node. The exit code (0 for success, non-zero for fail) signals success/failure to that orchestrator. That integration is straightforward because Elspeth is a CLI tool.
- If an experiment is extremely long (hours), one might run it on a screen or background process. Checkpointing means if the process or VM fails mid-run, one can start another process on same or different machine with the same config and the checkpoint file to continue. So physically, to resume on a different machine, you'd need the checkpoint file transferred to that machine and access to the same environment (like same data and output destinations).
- The system is not resource-elastic by itself; but if needed one could run multiple processes in parallel on a big VM to use more cores, etc. The controlling factor is memory for each.

- **Security of the Host:**

- The machine/container should be secured because it handles potentially sensitive data in memory and outputs. Standard OS hardening (patching, limited user access, etc.) should be in place. The machine should also be in a secure network segment if processing sensitive data.
- Keys and credentials should be provided securely (e.g., via environment variables set by a secure mechanism, not hardcoded in scripts). In containers, using something like Kubernetes secrets mounted as env vars is common. On VMs, using managed identity or storing secrets in Azure Key Vault and fetching them at runtime (outside Elspeth's scope, but possible by customizing how environment is prepared).
- If outputs need to be protected at rest, the physical environment should implement disk encryption (or one could use an encryption sink plugin, but none is standard, so relying on infrastructure encryption is prudent).
- Logging on the host: ensure logs do not capture sensitive content inadvertently (Elspeth's default logs are mostly high-level, but running with `--log-level DEBUG` might log prompt content or entire responses due to library debug logs, so that should be avoided on production runs with real data unless needed and then logs treated as sensitive).

• Physical Data Flow:

- Input data source (like a blob file) flows over network to the Elspeth host and is loaded into memory (DataFrame).
- Prompt content (from that data or from config prompt templates) flows from host to LLM service over network.
- LLM responses flow back to host and reside in memory, then in output files (CSV, etc.).
- Output files might then flow to external sinks (like uploaded to repo or blob) over network.
- Throughout, data in transit is encrypted (HTTPS) and data at rest on host can be encrypted by OS (if using encrypted disk).
- If using signing, the private key used to sign must reside on the host (maybe in memory or file) – physical security of that key is important. It should be on disk encrypted or in a secure store. Perhaps an HSM integration could be possible if signing plugin were extended (not currently done, but a consideration for physical key management).

• Resource Allocation:

- If running on a container platform, you would allocate CPU/memory to the Elspeth container. For instance, give it 4 CPUs (to match concurrency 4) and say 8GB RAM if processing moderately large data. If data is huge, scale accordingly.
- One should also monitor external API usage – e.g., if running many experiments, ensure the API provider (OpenAI etc.) supports that volume or concurrency from one IP (to avoid being throttled beyond what Elspeth handles).

Example Deployment Scenarios: 1. *Local Testing:* A data scientist runs Elspeth on a laptop. The config might point to a local CSV and use the OpenAI API with an API key loaded from `.env` file. The CLI writes outputs to a local `outputs/` folder. Here, the system boundary is just the laptop; internet is used for API calls. This scenario is simple but not likely accredited for sensitive data (maybe for dev). 2. *On-Prem Secure Server:* Elspeth is installed on a hardened RHEL server in a data center. The server has no direct internet; instead, it uses an HTTP proxy to reach OpenAI. The OpenAI key is stored in a root-owned file and loaded to env by a wrapper script when the experiment runs (so regular user running Elspeth doesn't see the key). The server is in a VLAN that can reach internal blob storage and the proxy for external. When Elspeth runs, it downloads data from an internal blob (no internet needed for that), calls OpenAI via the proxy (the requests library picks up `HTTPS_PROXY` env var), then outputs results

to an NFS share that is secured. If configured, it could also push results to an internal Git server (accessible without leaving intranet). This scenario fits many corporate setups. 3. *Azure Cloud VM or Container*: Elspeth runs on an Azure Virtual Machine that has a managed identity. The config uses Azure OpenAI (so no static API key needed, it uses the identity's token) and Azure Blob (again using identity). The VM has outbound internet open to Azure services. The experiment runs entirely within Azure, and outputs might be written to Azure Blob or Azure DevOps repo (via service connection). The VM's disk is encrypted by Azure. This could satisfy many government cloud requirements if properly configured (no secrets on disk, all auth via identity, TLS everywhere).

Physically, Elspeth's simplicity (no microservices, no always-on components) means it's quite straightforward to deploy and tear down. From an ATO perspective, fewer physical components often means fewer attack surfaces to secure: - There's no database or cache to secure. - No web server to harden. - Only the execution environment and network egress need strong controls.

Finally, regarding scaling physically: if needed to handle more load, one could horizontally scale by running multiple instances of Elspeth on separate machines or containers (like parallelizing different experiment suites on different nodes). Because there's no shared state, this scales linearly. If an ATO environment requires redundancy, one could have two servers and run identical experiments on both to compare and ensure determinism (just an idea – for verifying results). But typically, high availability is not crucial since one can re-run an experiment if a machine fails mid-run (thanks to checkpointing and idempotence of operations except time/cost).

In summary, the physical deployment of Elspeth is that of a **self-contained batch application**. Its interactions with external systems are through standard network channels (HTTPS), and it relies on the host environment for security (OS, container runtime) and scheduling (triggering runs). This straightforward physical architecture can be easily integrated into existing enterprise infrastructure and controlled with existing security measures (firewalls, proxies, OS-level encryption), making it feasible to obtain an ATO for use in sensitive environments.

2.3 Data Architecture (Information View)

Elspeth's data architecture revolves around the flow and transformation of data through the system – from input data sources to in-memory structures to output artifacts. The system primarily handles experiment configuration data, input datasets, prompt/response data, and result metrics. It does not maintain a long-term database; rather, it performs **extract-transform-load (ETL)** processes in memory and produces outputs to files or external sinks. Key aspects of the data architecture include:

Data Input and Models: - **Configuration Data:** The experiment settings (YAML profile) is a key data structure. It contains nested mappings for data source, prompts, LLM parameters, plugins, etc. This is parsed into a `Settings` dataclass which has fields like `datasource`, `llm`, `sinks` (list), `suite_defaults`, etc. ²³⁰. This structured config is then used to instantiate internal objects. The structure of config is validated by `validate_settings()` using rules defined in code (like all sink entries must have `security_level` as enforced in code ⁴⁶). Thus, configuration data flows from YAML text to an internal object graph connecting to plugin instances. The data model here is that each plugin definition in YAML is essentially a small record with fields (name, options, security_level, determinism_level) ^{296 298}. The loader transforms those into actual plugin objects and attaches context (e.g., `_elspeth_plugin_name`, `_elspeth_security_level` as attributes on plugin instances) ²⁸⁸. - **Input Dataset:** The primary input is typically **tabular data** – often a CSV or Excel file or a database table – containing rows on which experiments run. For example, a dataset might have columns like `id`, `user_query`, `expected_answer` (if evaluating model), etc. Elspeth uses pandas

to load these into a `pd.DataFrame`. This DataFrame (with shape [N rows x M columns]) is the main data structure for inputs in memory. It's attached to a `security_level` attribute (the code sets `df.attrs['security_level'] = level` if the data source had one) ¹⁴⁸, and possibly a `schema` attribute if provided by DataSource (some data sources might supply a JSON schema or data schema object in `df.attrs['schema']`). This DataFrame is passed to `ExperimentRunner.run()`, which then iterates over it. The architecture treats this dataset mostly as an opaque collection of records; operations like slicing, shuffling, etc., could be applied (e.g., orchestrator truncates to `max_rows` if set by simply doing `df = df.head(max_rows)` ¹³³). No relational database or advanced indexing is used beyond pandas capabilities.

- **Prompt Composition Data:** Within each row, fields may be combined into a **prompt context**. The `prepare_prompt_context` function creates a dictionary for each row, mapping field names to values, possibly filtering by `prompt_fields` list if provided ¹⁶³. So, for each row's data, a flat dict (or nested if input fields are nested JSON, but presumably flat keys) is the data passed into the prompt template. Also, any fixed prompt templates from config (system and user templates in `OrchestratorConfig.llm_prompt`) are stored as strings or compiled Jinja templates ^{191 194}. The `PromptEngine.compile()` produces a `PromptTemplate` object which contains: - The raw template text, - A Jinja2 template object, - A list of `required_fields` the template expects ^{5 461}, - Defaults provided for missing fields (if any), - Metadata (like a name). This compiled template is data that the runner holds (in `self._compiled_user_prompt` etc.) for rendering each row. This ensures template parsing is done once, and then filling (rendering) is done repeatedly. Data flows from the row context dict into the template to produce a final prompt string. Jinja2 uses the context mapping to substitute variables – an implicit data processing step (with `StrictUndefined` meaning if any required variable missing in context, it raises an error immediately) ^{462 463}. If an error (missing field) occurs, it's caught and turned into a `PromptRenderingError`, which runner catches and treats as a failure for that row ³⁶⁰.

- **LLM Request/Response Data:** When calling the LLM API, Elspeth constructs an `LLMRequest` object with fields: - `system_prompt` (string), - `user_prompt` (string), - `metadata` (a dict containing things like attempt number, maybe experiment name or other context) ^{141 166}. This is a transient data structure passed through middlewares. The integration with openai API converts it to an HTTP payload (usually JSON), which is sent out. The response from the LLM comes as either: - If using OpenAI SDK: a Python object (OpenAI returns a `Completion` or `ChatCompletion` object from which content, usage, etc., can be extracted). - If using raw requests: a JSON payload which plugin code interprets. In either case, Elspeth normalizes it into a dictionary called `response` with keys like `content` (the text output), maybe `finish_reason`, and importantly any `metrics` like token counts or model usage. The runner's `_execute_llm` code explicitly updates `response.setdefault("metrics", {})["attempts_used"] = attempt` and sets a `response["retry"]` dict with history ^{464 370}. So the response data structure that flows into further processing (plugins, output) is a dict containing at least: * `"content"`: the main model output (for chat, maybe combined or just the assistant content), * `"metrics"`: a dict of numeric metrics from the model (e.g., token counts, cost, perhaps latency if `cost_tracker` adds it), * `"retry"`: info about how many attempts it took and the timeline of attempts, * Possibly `"prompt"` if dummy LLM returns it (like `DummyLLM` returned prompt echo ⁴⁶⁵), * Possibly `"error"` if an error was caught and we are packaging it as failure (though in code they create separate `failure` dicts, not put `"error"` in response unless it's a recovered scenario). Each record in results ultimately is represented as a dictionary (for JSON serialization ease). The CLI's `_result_to_row` function flattens these for CSV output: - It pulls `record["response"]["content"]` into a column `llm_content`, - It pulls any named alternate responses into `llm_{name}` columns ¹⁸⁸, - It flattens any nested metrics into `metric_{name}` columns ¹⁷³, - It copies `record["security_level"]` to a column if present ⁴⁶⁶. So, internally, each result `record` is a nested dict, but for output to a table it becomes one flat row. For JSON output, they likely keep structure (the JSON output could mirror the internal dict structure, which is easier to interpret programmatically).

- **Failure Data:** If a row fails (due to exception in prompt rendering, LLM call exhaustion, validation fail, etc.), the runner creates a `failure` dict containing: -

"row": the input context (so one knows which input failed), - "error": the error message string, - "timestamp": when it failed, - Possibly "plugin" if early_stop reason includes plugin name, etc. Failures are collected in a list separate from results ³²⁴. They are included in output metadata (so the summary knows how many failures and can list them) ³⁴⁸. The SuiteReport may output a failure_analysis.json or markdown to detail them ⁴⁴⁸. This separation of normal results vs failures is part of data architecture – it ensures that one can easily see incomplete results versus successes. For CSV output, failures might not appear (since they have no normal result values; the CLI prints errors to log for failures). - **Aggregated Metrics:** After processing all rows, aggregator plugins produce summary metrics per experiment, which are stored in a dict `aggregates` in the payload ³⁴⁴. For example, an aggregator might produce `{"average_score": 0.82, "max_latency": 1.5}`. These become top-level fields in the output JSON under a "aggregates" key, and can be reported in summary documents. If baseline comparisons are done, the comparisons are stored as a nested dict under each experiment's result (payload["baseline_comparison"] with keys for each baseline plugin name) ¹⁰². E.g., `baseline_comparison: {"score_delta": 0.05}` meaning experiment's score is 5% higher than baseline. The data architecture thus includes hierarchical output: experiments -> metrics -> baseline comparisons, etc., which can be navigated in JSON or flattened for reports. - **Metadata and Audit Data:** The runner attaches a `metadata` dict to payload containing run-level info: - "row_count" (or named "rows" and "row_count" both set to number of results) ²⁰⁸, - "retry_summary" (with total requests, total retries, and exhausted count) ²¹, - If cost tracking is enabled and had data, "cost_summary" with aggregate cost/tokens ⁹³, - If early_stop triggered, "early_stop": {reason...} with plugin name and cause ³⁵⁰, - "security_level" and "determinism_level" of the run (the resolved values, often from config) ¹⁴⁸, - Possibly "failures" listing failure records if any ⁴⁶⁷. This metadata provides a comprehensive summary that goes into final output or logs. For instance, if writing an analytics report, these fields feed into an executive summary ("Processed N rows with M successes and K failures. Total tokens: X, cost: \$Y." etc.). - **Output Data Models:** - **CSV/Excel:** Tabular outputs have columns for key metrics and results. The CLI's flattening scheme defines those column names as described. For instance, output CSV might look like: | id | user_query | llm_content | metric_score | retry_attempts | security_level | and so on. This is convenient for analysts to open in Excel. If an experiment had multiple named responses (say it called two LLMs for each input), the CLI flattening puts each under `llm_modelA`, `llm_modelB` columns. - **JSON Output:** If using `--export-suite-config` or similar, they can export the internal configuration or results as JSON. We see SuiteReportGenerator likely producing `.json` files: * `consolidated/validation_results.json` (maybe a JSON of all validation warnings across experiments) ⁴⁶⁸, * `analytics_report.json` (structured metrics), * etc. So JSON is used as an interchange format for capturing the nested output data without loss (unlike CSV which has to flatten). - **Markdown/HTML Reports:** Analytical sinks might produce human-readable documents summarizing runs (ex: `executive_summary.md`, `comparative_analysis.md`). These are derived from the data in metadata and aggregates. E.g., `executive_summary` might list each experiment's key metrics and baseline deltas. Visual HTML might embed base64 images or link to PNGs produced. - **Signature Artifact:** If signing is used, the data artifact might be a signature file (like a `.sig` or `.json` manifest containing hashes). For instance, mermaid shows: Signing uses HMAC manifest to produce "Audit Artefacts" ⁴⁵⁰. Possibly the Signed sink writes a JSON manifest of all output files with their SHA-256 hashes and an HMAC using a secret, and that manifest is an audit artifact. This manifest's data model would be something like:

```
{
  "artifacts": [
    {"name": "results.csv", "sha256": "<hash>"},
    {"name": "analysis.xlsx", "sha256": "<hash>"}
  ],
  "generated_at": "...",
```

```
"hmac_signature": "<hmac over all above data>"
}
```

This would allow later verification of any output file by recomputing its hash and checking against manifest and verifying HMAC with the secret. This is a data artifact solely for integrity verification (the "Audit evidence"). - **Traceability Matrix:** There's mention of `TRACEABILITY_MATRIX.md` and `CONTROL_INVENTORY.md` in docs (perhaps prepared manually or semi-automatically). Those are documents mapping controls to code and evidence. If Elspeth produces any data for those (maybe not automatically, likely static docs), it's more a part of system documentation rather than run-time data.

Data Flow & Transformation: - **Extraction:** Data is extracted from sources (file, blob, etc.) via DataSource plugin and turned into DataFrame in memory ⁸⁵. - **Transformation:** - Each data row is transformed into an LLM input (prompt) by combining it with the prompt template. That is a transformation from structured row to text prompt string (with some logic like filling fields, maybe formatting numbers, etc., via Jinja2 filters). - The LLM response is transformed into structured result: splitting content vs. metrics, tagging with security level, etc. Also merging any extra computed metrics (like if row plugin returns `{"sentiment": "positive"}`), that is merged into the result record's "metrics" sub-dictionary). - Summaries (aggregator output) transform a list of metrics into a single aggregate metric. - Baseline comparison transforms two sets of aggregates into differences or statistical measures. - Sanitization is a transformation applied to output data (e.g., the CSV sink scanning each cell of output and prefixing dangerous characters) ¹¹⁰. This ensures the output data has no malicious formula - it's a data transformation for security. - The signing sink transforms a set of output files into a cryptographic manifest (the transformation is from arbitrary file bytes to their hash and then to a signature). - The zip sink transforms a set of files into a combined archive. These transformations are all deterministic (or labeled if not; but by default, Elspeth aims for deterministic outputs given same inputs and config, aside from nondeterminism in LLM if allowed). - **Load (Outputs):** Data is "loaded" into output sinks, which means writing to files or uploading via API. For example: - CSV sink loads data into a CSV file on disk. - Excel sink loads data into an .xlsx file (with multiple sheets possibly). - Repo sink loads data into a remote repository (committing a JSON or CSV into a repo path). - Visual sink loads images into .png files, and an HTML summary linking them. - Signed sink loads a manifest and signature into a file. All these final forms are the "persisted" data after Elspeth's run. They are the artifacts that stakeholders ultimately use or verify.

Data Volume and Performance: - The design leverages Pandas for potentially large data sets; however, extremely large (millions of rows) could be slow or memory-heavy, which is acknowledged by their mention of future streaming architecture (not present yet) ⁴⁶⁹. For now, all data is in memory. The checkpoint allows chunk processing in multiple runs if needed (one could artificially split input into chunks and run sequentially, using checkpoint to skip done parts, as a workaround to not loading all at once). - The data processed by LLM is typically much smaller (prompts are usually not more than a few thousand characters, and responses similar). The heavy part is that it might process N of those where N is dataset size. So the volume to LLM API = N * average prompt+response size. This can be large but is handled via iterative calls and not stored all at once except in aggregated form. - The results dataset has as many records as input (minus failures). So similar scale to input. If input was 10k rows, output will be 10k result entries (which can be saved in a CSV of comparable size to input CSV plus extra columns). - If aggregator and baseline produce small additional data, negligible overhead. - Possibly the largest output might be logs or traces if an LLM returns very large text or if debug logging is on. But those are typically within expected size (maybe tens of KB per response at most). - As for data integrity, each transformation tries to preserve relevant information: - StrictUndefined ensures no missing field goes unnoticed (so data required by prompts must be present). - Validation plugins ensure that if output data doesn't meet certain structure (like JSON parseable or contains required keys), it is flagged. This means

the output data is *integrity-checked* at a semantic level. - At the end, the signature ensures the integrity at storage level (bits not tampered). - The classification tagging ensures all data in outputs is marked by sensitivity level, which is then used by artifact pipeline to prevent copying to lower security sinks (so data confidentiality and proper handling is enforced through data labeling).

Traceability: - Elspeth's data architecture was clearly designed for traceability: for each output row, one can trace back to input (`record["row"]` stores input fields) ⁴⁷⁰, and see what happened (metrics and potentially error info). The presence of `history` in retry, and `failures` list capturing problematic inputs, ensures nothing is lost silently. - The metadata with cost and attempts gives overall trace of resource usage.

In conclusion, Elspeth's data architecture is about converting **structured inputs to structured outputs with an LLM in the loop**, while enriching that data with metrics and enforcing constraints. It flows as: Config (static data) + Input (dynamic data) -> (via templates) prompt text -> LLM output text -> structured results (with metrics) -> aggregated metrics -> outputs (files/records) and audit info. Each step is captured in data structures that are either in-memory (during processing) or persisted (in outputs and logs). There is no long-term state beyond outputs, making the data architecture relatively simple in terms of data lifecycle: data comes in, is transformed, and goes out, without lingering in the system outside of runtime memory and designated output artifacts.

3. Component Catalog

In this section, we provide specifications for each major component identified in the architecture, describing their purpose, responsibilities, interfaces, and interactions, with references to the source code for evidence.

- **CLI Interface** (`elspeth.cli`) - *Purpose:* Serve as the command-line entry point for users to run experiments or suites. *Responsibilities:*
 - Parse user arguments (using `argparse`), such as `--settings` (path to YAML config), `--profile` (which profile in config), `--reports-dir` (where to output reports), `--disable-metrics` (flag to modify settings), etc. ²⁴ ⁶⁹.
 - Invoke configuration loading and orchestrator execution. The CLI is essentially a thin wrapper that calls `validate_settings()` and `load_settings()` to get a `Settings` object, then decides to run a suite or single experiment based on presence of suite config. ²⁶ ²⁵.
 - Handle top-level control flow: e.g., if `--validate-schemas` is provided, it calls a function to perform schema validation on the input and exits without running experiments (this is a dry-run mode to verify input compatibility) ²²³. If `--create-experiment-template` is given, it delegates to a utility to add a template experiment in the suite config (then exits) ²⁵². These ancillary features help manage experiment definitions and verify data without executing everything.
 - Configure logging verbosity according to `--log-level` argument ⁴³³. For example, `--log-level DEBUG` makes internal debug logs visible (useful for troubleshooting).
 - Manage output directory setup: If `--reports-dir` is specified, it ensures the directory exists, and passes it along to `SuiteReportGenerator` for writing reports ⁴⁷¹.
 - After orchestrator run, handle final output tasks: The CLI might print a snippet of results for quick view. In fact, it prints a preview of the output DataFrame (first few rows) if `--head N` is not 0 ⁴⁷². This gives the user immediate feedback in the console.

- Exit with appropriate status code: if any errors occurred during run, the CLI will propagate the exception or exit with non-zero (the code calls orchestrator and doesn't catch exceptions like `SystemExit` from config validation or runtime errors, so the process will exit accordingly).

Dependencies/Interfaces: The CLI uses `argparse` interface to gather args. It relies on `elspeth.config` for settings loading ⁴⁷³, and on `ExperimentOrchestrator` / `ExperimentSuiteRunner` for executing runs. It calls `ExperimentSuiteRunner.run()` or `ExperimentOrchestrator.run()` depending on config (if suite defined) ²⁴⁸ ⁴⁰. It also directly interfaces with the filesystem (for reading YAML, writing any quick outputs). *Interactions:* The CLI doesn't directly manipulate experiment internals beyond initial setup; it delegates those to deeper layers. For example, after run, if multiple experiments and `--reports-dir` given, it invokes `SuiteReportGenerator.generate_all_reports()` to produce consolidated reports ³³, indicating that CLI passes control to a reporting utility (which is logically a sink).

- **Configuration Loader** (`elspeth.config.load_settings`) - *Purpose:* Read the YAML configuration and construct the internal `Settings` object with all necessary runtime objects for an experiment suite. *Responsibilities:*

- Parse YAML using `yaml.safe_load`, then select the relevant profile section (default or as specified) ⁴²⁷ ²⁵⁴. The config is expected to be in profiles (mapping of profile name to settings). The loader picks the requested profile's dict.
- Validate the structure and combine parts of config:
 - It uses `validate_settings()` (which internally uses `elspeth.core.validation.validate_settings` logic) to produce warnings or errors for any config issues (like missing required fields) ²²². For instance, if a plugin definition lacks `security_level`, `validate_settings` will raise an error, preventing loading until fixed.
 - It collects **Prompt Configuration** (prompts, prompt_fields, criteria) via `_collect_prompt_configuration` ⁴⁷⁴ and **Plugin Definitions** (row, aggregator, etc. plus concurrency, early_stop, rate_limiter definitions) via `_collect_plugin_definitions` ²⁵⁸ ⁴⁷⁵. These functions extract those sub-dicts from the profile data for further processing.
- Instantiate plugins and controls:
 - It calls `_instantiate_plugin(def, "datasource", datasource_registry.create)` for the datasource definition ¹³⁰. Similarly for `llm` definition with `llm_registry.create` ¹³⁰, and for each sink in `sink_defs` via `_instantiate_sinks()` which uses `sink_registry.create` for each ⁴⁷⁶ ¹²⁴.
 - It calls `create_rate_limiter(rate_limiter_def)` and `create_cost_tracker(cost_tracker_def)` if such definitions exist in config ⁶⁷, to instantiate those control objects (which might be classes with state, e.g., a `TokenBucket`).
- Build the `OrchestratorConfig` dataclass (which is embedded in `Settings`) with all resolved fields: The code returns an `OrchestratorConfig(...)` populated with `llm_prompt` templates (system & user strings), `prompt_fields` list, `criteria` list, lists of plugin definitions, concurrency config dict, `early_stop` config, etc. ²⁶² ²⁶³. This is essentially a normalized config ready for orchestrator use.
- Prepare suite-level constructs:
 - If `suite_defaults` present in config, merges them with any selected prompt pack (the config may reference a `prompt_pack` name that indexes into a `prompt_packs` dictionary in config). The loader uses `_prepare_suite_defaults` to overlay a chosen prompt pack's values onto the base `suite_defaults` ⁴⁷⁷ ²⁶⁰. This results in a `suite_defaults` dict inside `Settings`.

- It extracts `suite_root` path if given (for directory of suite experiments files) ⁴⁷⁸ and ensures it's a `Path` object.
- Ultimately, it returns a `Settings` object (a dataclass defined in `elispeth.config.Settings` with fields like `datasource`, `llm`, `sinks`, `orchestrator_config`, etc.) ²³⁰ ¹³¹. The `Settings` encapsulates everything the CLI and Orchestrator need to proceed. *Dependencies:* PyYAML for parsing YAML text ²⁶⁴. The plugin registries (which are singletons or module-level factories defined in core) for creating plugin instances ²⁶⁵ ²⁶⁶. The validation module for config correctness checks. The dataclass definitions (`Settings`, `PromptConfiguration`, `PluginDefinitions`) for type structure ²⁵⁵ ²⁵⁶. *Interactions:* The Config Loader is invoked by the CLI (or by Orchestrator in some tests) to supply initial state. It interacts with no external services (except reading local config file from disk). It heavily interacts with plugin layer: for each plugin listed in config, it calls the appropriate factory and may catch errors if a plugin name is unknown (in which case registry would likely throw and loader then raises `ConfigurationError`). It sets security and determinism defaults for plugins if not provided by config via helper `_prepare_plugin_definition` (which merges any `security_level` from definition and options and normalizes them, raising error if invalid) ⁴⁷⁹ ⁴⁸⁰. *Example:* If config has:

```
datasource:
  plugin: "csv_local"
  options: {path: "data/input.csv"}
llm:
  plugin: "openai"
  options: {model: "gpt-3.5-turbo", api_key: "${OPENAI_API_KEY}"}
sinks:
  - plugin: "csv_file"
    options: {path: "outputs/results.csv"}
    security_level: "OFFICIAL"
    determinism_level: "guaranteed"
```

The loader will create:

- A `CSVLocalDataSource` instance via registry ¹³⁰.
- An `OpenAIClient` instance via `llm_registry` ¹³⁰.
- A `CsvFileSink` instance via `sink_registry` ¹²⁴, applying the given security/determinism (or raising if they were missing). It then populates `Settings` with those plus default `cost_tracker` and `rate_limiter` (if not in config, those remain `None`). This `Settings` is returned to CLI which then passes it to Orchestrator.
- **ExperimentOrchestrator** (`elispeth.core.orchestrator.ExperimentOrchestrator`) - *Purpose:* Coordinate execution of a single experiment (one dataset through one LLM pipeline and associated plugins). *Responsibilities:*
 - Initialize experiment execution context and dependencies:
 - Stores references to the `datasource`, `llm_client`, list of `sinks`, `rate_limiter`, `cost_tracker` provided via `Settings` ²⁷³ ²⁷².
 - Resolve the combined security level of the experiment from its components. It calls `resolve_security_level(datasource_sec, llm_sec)` to determine an overall classification for experiment data ⁴⁸. This picks the most restrictive (e.g., if `datasource` or `llm` has higher level, that wins). It then creates an `PluginContext` object capturing

- experiment name, type "experiment", and that security level, plus provenance (like `orchestrator:experimentName.resolved`)¹²⁵. This context is essentially metadata passed to plugins for logging and enforcement. The Orchestrator attaches this context to the `rate_limiter` and `cost_tracker` by calling `apply_plugin_context` on each (so that these control objects know which experiment they belong to, which can help in logging or separate accounting)^{75 275}.
- Instantiate internal plugins (row, agg, val, early_stop) as per OrchestratorConfig:
 - It iterates over each `row_plugin_def` in config and calls `create_row_plugin(defn, parent_context=experiment_context)` from the plugin_registry⁷³. Each returns a plugin instance (with possibly its own internal state). Similarly for `aggregator_plugin_defs`, `validation_plugin_defs`, and `early_stop_plugin_defs`^{481 276}. The resulting lists of plugin instances are stored in local variables.
 - It also creates any LLM middleware by calling `create_middlewarees(llm_middleware_defs, parent_context=experiment_context)`, which returns a list of middleware instances (some may be shared if definitions repeat – the factory might cache, but orchestrator likely just gets fresh ones unless plugin registry ensures singletons for certain global ones)¹⁰³.
 - It now has lists: `row_plugins`, `aggregator_plugins`, `validation_plugins`, `early_stop_plugins`, plus the single `rate_limiter` and `cost_tracker` objects (with context applied).
 - It constructs the **ExperimentRunner** (or uses one passed for testing) with all these components:
 - Passes `llm_client`, `sinks`, `prompt_system`, `prompt_template`, `prompt_fields`, `criteria`, `row_plugins`, `aggregator_plugins`, `validation_plugins`, `rate_limiter`, `cost_tracker`, `experiment_name`, `retry_config`, `checkpoint_config`, `llm_middlewarees` (the list from earlier), `prompt_defaults` (like default values for template), `concurrency_config`, `security_level`, `early_stop_plugins`, `early_stop_config`^{74 119}.
 - This call to `ExperimentRunner(...)` returns a runner instance, which the orchestrator then attaches the `plugin_context` to as an attribute (for potential use by runner or plugins)²⁷⁷.
- Execute the experiment run:
 - Calls `df = self.datasource.load()` to retrieve the input DataFrame⁸⁵. Catches exceptions from DataSource if any (e.g., file not found – it would raise an error; orchestrator doesn't explicitly catch here, so if DataSource raises, it propagates and aborts run).
 - If `self.config.max_rows` is set (OrchestratorConfig may contain a `max_rows` limit from YAML), it truncates the DataFrame to that number of rows (using pandas `.head(max_rows)`)¹³³. This prevents processing more data than desired.
 - Then, if `df.attrs['schema']` exists and `--validate-schemas` was requested (the CLI actually handles schema validation separately via a flag, but orchestrator itself doesn't do it here unless in schema-only mode), orchestrator could validate plugin compatibility, but actually `validate_settings` does a lot of validation earlier. The code shown doesn't explicitly call schema validation here (instead, in CLI's validate mode it calls a separate function).
 - Calls `payload = self.experiment_runner.run(df)`¹³⁴. This is the core hand-off to the runner which processes the data fully and returns results and metadata as

described in component 3. After this returns, orchestrator does not further transform data (the runner returns final payload).

- It then returns that `payload` (or nothing and handles output itself? Actually, Orchestrator's `run()` returns the payload dict up to CLI or suite logic) ⁴⁸².
- Error handling: The orchestrator doesn't catch exceptions from runner explicitly; if runner raises (due to unhandled plugin error, etc.), it will propagate up. Typically, runner catches most plugin errors and returns them in payload rather than raising. Orchestrator, however, might catch the specific case of schema validation mode and just return after doing it, but that's in CLI logic, not orchestrator (or orchestrator might be used in that context too). *Dependencies/Interfaces*: It heavily uses plugin registry factories (via calls like `create_row_plugin`) ⁷³. It interfaces with `PluginContext` from `elasticsearch.plugins.PluginContext` to propagate context info ¹²⁵. It relies on the `ExperimentRunner` interface: specifically, the runner must implement a `run(dataframe)` method returning a payload dict ¹³⁴. It interacts with DataSource object by calling its `load()`, and expects a pandas DataFrame in return ⁸⁵. It interacts with RateLimiter and CostTracker only to apply context (and later runner will use them; orchestrator doesn't actively throttle, it just passes them in). *Interactions*:

- With DataSource: orchestrator triggers data extraction from external data.
- With Runner: orchestrator delegates actual data processing to it.
- With RateLimiter/CostTracker: orchestrator sets them up but does not drive their behavior beyond initialization.
- With CLI/Suite: If part of a suite, the ExperimentOrchestrator might be used implicitly via SuiteRunner building multiple orchestrators or directly multiple runners. In code, SuiteRunner doesn't actually create Orchestrator objects; it merges config and directly builds runner for each experiment, thus somewhat bypassing making a full Orchestrator instance for each (to avoid re-loading data for baseline, etc., it loads data once and uses it for all experiments? Actually, from suite_runner code, it loads suite baseline and experiments separately via ExperimentConfig – likely each experiment can have separate data if needed, but typically they share same input).

- **ExperimentSuiteRunner** (`elasticsearch.experiments.ExperimentSuiteRunner`) – *Purpose*: Manage execution of a group of experiments (suite), including baseline handling and consolidated reporting. *Responsibilities*:

- Interpret the `ExperimentSuite` object (built from config) which contains:
 - Possibly a baseline experiment definition (`suite.baseline`),
 - A list of other experiment definitions (`suite.experiments` list) ⁴⁸³.
 - These are `ExperimentConfig` objects capturing each experiment's specific config (prompts, LLM choices, etc.).
- For each experiment in the suite, build and run an ExperimentRunner:
 - Merges suite-wide defaults, any prompt pack values, and the experiment's own settings via `ConfigMerger`. This ensures each experiment's config is a combination of common defaults and specific overrides ²⁸⁰ ¹²⁹.
 - Creates experiment context similar to orchestrator (with security level resolved among experiment's config, pack's level, defaults) ²⁸⁴ ²⁸⁵.
 - Instantiates plugins for that experiment (row, agg, etc.) with context (calls `create_row_plugin` etc. within `build_runner()`) ⁷⁹.
 - Instantiates or reuses controls: The SuiteRunner is designed to either create a new RateLimiter/CostTracker if experiment-specific config is given, or reuse the one from suite defaults if not (the code shows if no new def and suite default has a shared object, apply

- context to it and use it). This way, experiments can share a global limiter or have separate if configured.
- Builds an ExperimentRunner for each experiment with those specific components ⁷⁹₂₉₄, same as orchestrator would, then returns that runner.
 - Execute experiments in sequence:
 - If baseline is defined, run it first (they arrange experiments list such that baseline is first) ³⁰³₃₁₄.
 - For each experiment config, including baseline and others:
 - Determine appropriate sink instances: The suite runner logic can override sinks per experiment. The code indicates it calls `_instantiate_sinks()` for each experiment definition if it has custom sink_defs; if not, it clones base sinks or uses a provided `sink_factory` lambda (from CLI) to copy base sinks for each experiment so outputs don't collide ⁴⁸⁴. E.g., CLI passes a `sink_factory` that when given an Experiment, renames the base CSV file to prefix experiment name ⁴⁸⁵.
 - After preparing sinks and runner via `build_runner()`, it attaches shared middlewares if needed (the code shows it maintains `_shared_middlewares` dict to reuse identical middleware across experiments, e.g., an Azure telemetry middleware might be shared so that it logs suite-level info only once) ⁸²₁₉₀.
 - It then runs `runner.run(df)` for each experiment (where `df` might be the same DataFrame for all if suite uses one dataset, or each experiment could have its own data in config; typically they share input data unless experiment config specifies a different datasource).
 - Collect the returned `payload` in a results dictionary keyed by experiment name ³¹³₃₁₅.
 - Call `on_experiment_start` and `on_experiment_complete` hooks on any middleware that implements them (like logging a start event with experiment metadata, and logging completion with payload) ¹⁰⁷₁₀₅.
 - After all experiments, perform baseline comparisons:
 - If baseline exists, for each non-baseline experiment, it iterates through any baseline comparison plugin definitions, creates plugin instances (via `create_baseline_plugin`) and calls `plugin.compare(baseline_payload, experiment_payload)` ⁸⁰₁₀₁. Each plugin returns a diff (e.g., dict of metrics differences). The suite runner then stores that in the experiment's result under `"baseline_comparison"` key ¹⁰¹₃₁₆.
 - It also calls any middleware hook `on_baseline_comparison` after computing these, in case middleware need to log or adjust something after comparisons ⁴⁸⁶.
 - Gather suite-wide metadata:
 - The suite runner may compile a `preflight_info` at start (with count of experiments, baseline name) and call `on_suite_loaded` on shared middlewares (so they know about entire suite context) ⁴⁸⁷₃₁₀.
 - After experiments, it calls `on_suite_complete` on any middleware that has that (letting them finalize any suite-level logging or resource cleanup) ⁴⁸⁸.
 - It then returns the `results` dictionary containing each experiment's payload (with their results, metrics, comparisons, etc.) ⁴⁸⁹.
 - The SuiteRunner does not directly handle output sinks beyond instantiating them for experiments. The actual writing of artifacts was done in each experiment's runner already. However, a top-level reporting might be triggered by CLI after obtaining this results map (e.g., SuiteReportGenerator uses the results map to produce combined reports). *Dependencies:* Uses `ConfigMerger` (a helper that merges config layers, likely handling precedence of suite defaults vs experiment vs pack) ²⁸⁰₄₂₁. It relies on plugin registries to create each experiment's

plugins, and on `ExperimentRunner` for actual running. It uses `PluginContext` for each experiment as well to provide context to plugins (constructed similarly to Orchestrator's usage).
Interactions:

- With Orchestrator: Actually, SuiteRunner logically replaces Orchestrator when running multiple experiments. It doesn't call Orchestrator; instead, it replicates some of Orchestrator's steps for each experiment (like plugin instantiation, runner creation).
 - With Data Sources: If all experiments share one input dataset, ideally it should load it once and reuse. In current code, the `ExperimentSuiteRunner.run(df, defaults)` signature suggests you can pass in a DataFrame to use (the CLI does so: it calls `suite_runner.run(df, ...)` with df loaded by earlier in CLI's `_run_suite` ^{40 397}. Indeed, CLI loads the data once via `settings.datasource` before calling `_run_suite`, meaning SuiteRunner doesn't call `datasource` itself; it receives df from CLI. This is an optimization to not reload input per experiment. So interactions: CLI loads `df` via settings (thus Orchestrator logic for data load is bypassed in suite mode; they treat input as static across experiments unless experiment config had its own `datasource`).
 - With Rate Limiter / Cost Tracker: It may share these across experiments. E.g., if suite defaults had a `RateLimiter` object, it applies it to each experiment (so that if experiments run sequentially, the limiter continues counting usage across them if configured). If each experiment has separate def, it creates separate ones. Similarly for cost.
 - With output sinks: It ensures each experiment's outputs don't collide by using unique sink instances or calling the provided `sink_factory` (like renaming output files per experiment) ¹¹⁶. It doesn't itself combine outputs into one, except via separate reporting sinks (like a `VisualAnalytics` sink may be configured to aggregate metrics from all experiments; in design, a specialized plugin or the `SuiteReportGenerator` does that after run).
- **ExperimentRunner** (`elasticsearch.core.experiments.ExperimentRunner`) – *Purpose:* Execute the core loop of processing input records through LLM and plugins, managing concurrency, retries, and assembling results. *Responsibilities:*
- Initialize run state:
 - If `checkpoint_config` is provided (with a path and an ID field), load the set of already processed IDs from the checkpoint file (each stored line by line) ^{331 13}. Use this to skip those records.
 - Prepare data structures for concurrency and early stopping:
 - Set up `_active_early_stop_plugins` and `_early_stop_event` if any early_stop plugins configured ^{490 491}. It resets any prior run's state (in case same runner reused, though typically not reused across runs).
 - Initialize a thread-safe Event `_early_stop_event` and Lock `_early_stop_lock` to coordinate early stop flags across threads ^{329 492}.
 - Compile prompt templates using `PromptEngine`:
 - Use `PromptEngine.compile(system_prompt)` and `.compile(user_prompt)` to get compiled templates, stored in `self._compiled_system_prompt` and `_compiled_user_prompt` ^{194 493}.
 - For each entry in `criteria` (additional prompt variations), compile those templates too and store in `self._compiled_criteria_prompts` dict (keyed by criterion name) ¹⁶⁴.
 - These compiled templates also record required fields for validation (the engine uses `StrictUndefined`, but they also capture `required_fields` list to possibly check if defaults cover them – in code they consider raising if required fields not covered by defaults, but they put a TODO comment and do not fully enforce; still, it's info).
 - If input DataFrame has a `.attrs['schema']`, validate plugin compatibility with it:

- They call `self._validate_plugin_schemas(datasource_schema)` for each plugin type if schema is provided ⁹⁶. `validate_schema_compatibility` (imported as function) likely checks each plugin if it has a declared schema for input it expects (like a row plugin might say "I need fields X, Y of certain types"). If schema doesn't meet these, it raises or logs warnings. This ensures data is structurally compatible with plugins. The runner catches exceptions from it and would abort if failed (ensuring input schema issues cause early failure to avoid incorrect processing).
- Initialize structures for results:
- `rows_to_process` list to gather (index, row_series, context_dict, row_id) for each record that should be processed ^{163 494}.
- `records_with_index` list to collect output records (with original index).
- `failures` list to collect failure info dictionaries.
- Iterate through DataFrame rows (using `df.iterrows()` or `enumerate(df.itertuples())` - code uses `enumerate(df.iterrows())` to get index and row) ¹⁶³:
- For each row:
 - Construct `context` dict via `prepare_prompt_context(row, include_fields=self.prompt_fields)` which picks the columns relevant to prompt (if `prompt_fields` is None, likely include all) ¹⁶³.
 - Compute `row_id` if checkpointing is on: `row_id = context.get(checkpoint_field)` (e.g., if `checkpoint_field` is "id", it gets that value) ¹⁸³.
 - If `processed_ids` set is not None and contains this `row_id`, skip adding this row (already done in a previous run) ¹⁸⁴.
 - If `_early_stop_event` is set (some plugin triggered stop earlier), break out of the loop entirely (stop collecting new tasks) ³³³.
 - Otherwise, append `(idx, row_series, context, row_id)` to `rows_to_process` list ⁴⁹⁵.
- By end, `rows_to_process` is list of tasks we actually need to do (accounting for skipping and possibly early termination flag).
- Use concurrency configuration to decide how to execute:
- Check `if rows_to_process and self._should_run_parallel(concurrency_cfg, len(rows_to_process))` ¹³⁷. `_should_run_parallel` returns True if concurrency enabled and `max_workers > 1` and backlog \geq threshold ⁷.
- If parallel:
 - Call `self._run_parallel(rows_to_process, engine, system_template, user_template, criteria_templates, row_plugins, handle_success, handle_failure, concurrency_cfg)` ³³⁴.
 - This function sets up a `ThreadPoolExecutor(max_workers)` ^{335 336}.
 - Define `worker(data_tuple)` inner function which:
 - Checks `_early_stop_event` at start, returns if set to avoid wasted work ³³⁹.
 - Unpacks (idx, row, context, row_id) and calls `record, failure = self._process_single_row(...)` to actually process that one row ²³⁶.
 - Acquires a lock with `lock:` and then calls `handle_success(idx, record, row_id)` if record exists or `handle_failure(failure)` if not ³³⁷.
 - Submit all tasks to executor in a loop ¹⁵⁴, but before each submission:

- If `self.rate_limiter` and `pause_threshold > 0`:
 - It loops `while self.rate_limiter.utilization() >= pause_threshold: time.sleep(pause_interval)` to throttle submission speed ³³⁸. This effectively paces task dispatch to avoid overload on external API (utilization presumably returns fraction of allowed concurrent or allowed tokens used).
 - If `_early_stop_event` is set while scheduling, break out of submission loop (stop scheduling new tasks if stop triggered mid-schedule) ²⁴².
 - Once tasks submitted, the context manager closes (executor waits for tasks to finish because out-of-context).
 - If not parallel or no rows:
 - Use sequential loop: `for idx, row, context, row_id in rows_to_process: do:`
 - Check `_early_stop_event` before each iteration, break if set (stop processing further if stop triggered) ⁴⁹⁶.
 - Call `record, failure = self._process_single_row(..., context, row, row_id)` synchronously ⁴⁹⁷.
 - If `record` returned (not None), call `handle_success(idx, record, row_id)`; if `failure` returned, call `handle_failure(failure)` ³⁴³.
 - The `handle_success` closure does:
 - Append `(idx, record)` to `records_with_index` list.
 - If checkpointing: add `row_id` to `processed_ids` set and append `row_id` to checkpoint file (opening file in append mode and writing `row_id + newline`) ⁴¹⁰.
 - Call `self._maybe_trigger_early_stop(record, row_index=idx)` to check each `early_stop` plugin on this new record's data, possibly setting `_early_stop_reason` and `_early_stop_event` if any triggers ⁴⁹⁸. `_maybe_trigger_early_stop` locks around evaluating each plugin's check and logs if triggered, then sets `_early_stop_event` ^{499 500}.
 - The `handle_failure` closure does:
 - Append `failure` dict to `failures` list.
 - (It does not write to checkpoint, because failure typically means record wasn't processed to completion; we might consider not writing it to checkpoint so that re-run could attempt again).
 - If `failure` has a retry history, that's within `failure` and will be later summarized.
 - At end of parallel or sequential, all processed records (successful ones) are in `records_with_index`, and all failures in `failures`.
- Post-processing after record loop:
- Sort `records_with_index` by index and extract the `record` parts to a list `results` ⁵⁰¹. This ensures results are in original input order (since concurrency could have completed out of order).
 - Construct `payload = {"results": results}`. If failures list is not empty, add `"failures": failures` to payload ³²⁴.
 - Compute aggregates: For each plugin in `self.aggregator_plugins`, call `derived = plugin.finalize(results)` and if returned non-empty, put it into an `aggregates` dict under plugin's name ⁹². Add this `aggregates` dict to payload if not empty.

- Build `metadata` :
- `metadata["rows"] = len(results)` and `metadata["row_count"] = len(results)` ²⁰⁸ .
- Compute a `retry_summary` : iterate over every record in results and every failure, accumulate:
 - `total_requests = len(results) + len(failures)` (since each record had at least 1 attempt, each failure is an attempt that ultimately failed).
 - `total_retries = sum(max(info["attempts"]-1,0) for info in each record["retry"] if present + each failure["retry"] if present)` ²¹ . This adds up all extra attempts beyond first.
 - `exhausted = len(failures)` (requests that exhausted attempts and still failed).
- If any retries occurred (`retry_present` flag), set `metadata["retry_summary"] = {...}` with those totals ⁵⁰² ⁵⁰³ .
- If `aggregates` exists, also include it in metadata for convenience (so sinks can find aggregate metrics easily) ⁹³ .
- If `cost_tracker` exists, call `summary = self.cost_tracker.summary()` and if non-empty, attach `payload["cost_summary"] = summary` and `metadata["cost_summary"] = summary` ³⁴⁷ .
- If failures list is not empty, attach `metadata["failures"] = failures` (some redundancy with payload failures, but likely omitted to avoid duplication; the code shows they attach it though ³⁴⁸).
- Determine final security and determinism level of outputs:
 - Check `df.attrs.get("security_level")` (source data's level) and `self.security_level` (from config context). Compute `self._active_security_level = resolve_security_level(self.security_level, df_security_level)` ¹⁴⁸ . This ensures if data itself was labeled higher, that's respected. Set `metadata["security_level"] = self._active_security_level` .
 - Similarly for `determinism_level` : if data had an attribute or global config had one, combine them (the code sets `_active_determinism_level` similarly) and put in metadata ⁵⁰⁴ .
- If `_early_stop_reason` is not None (i.e., early stop triggered), attach `metadata["early_stop"] = copy of self._early_stop_reason` and also `payload["early_stop"] = reason` for completeness ³⁵⁰ .
- Attach `metadata` to payload as `payload["metadata"] = metadata` ⁵⁰⁵ .
- Prepare for output sinks:
- Build `bindings = self._build_sink_bindings()` which returns a list of SinkBinding objects, each containing sink id, plugin name, sink instance, artifact config (from config or introspection), and security_level for each sink ³⁵² ¹¹⁵ . If any sink has missing security_level, this will raise (ensuring all sinks declare it) ³⁹⁹ .
- Initialize `pipeline = ArtifactPipeline(bindings)` which computes execution order (resolving produces/consumes) ⁵⁰⁶ and store internal structure.
- Call `pipeline.execute(payload, metadata)` . This:
 - For each binding in resolved order, calls the sink's `write(results, metadata=metadata)` method ¹¹⁷ . It may filter which artifacts to pass based on `binding.consumes` (the pipeline likely passes either the whole payload or a map of requested artifacts; the code not fully shown, but logically, it finds artifacts in `ArtifactStore` for each sink's consume tokens and provides them).

- The `ArtifactPipeline` records any artifact produced by a sink using `ArtifactStore.register()` with proper id and security level ⁴⁴¹ ³⁹⁴. It then resolves these for later sinks requests via `resolve_requests()` ⁴⁴².
- It enforces security at production and consumption: For each dependency between sinks, it checks `is_security_level_allowed(producer.security_level, consumer.security_level)` when linking them ³, and `is_security_level_allowed(artifact.security_level, binding_clearance)` when passing artifacts to a sink ³⁸⁹, raising `PermissionError` if any violation.
- After `pipeline.execute`, all sinks have run. If any sink raised an exception (like a network error on push), it propagates (since `ExperimentRunner` doesn't catch around pipeline call) causing `run()` to raise – which would then bubble up to orchestrator and CLI.
- Finally, set `self._active_security_level = None` (clearing to avoid reuse confusion) ³⁵⁴.
- Return `payload` dict to caller (Orchestrator or SuiteRunner) ³⁵⁵. *Dependencies:* The runner uses many imported modules:
- `concurrent.futures.ThreadPoolExecutor` for concurrency ¹⁵⁴.
- `time.sleep` for throttling ⁵⁰⁷.
- Logging via `logger` for exceptions or early stop info ²⁰⁵.
- `json` not directly, but sinks might for output, not in runner code.
- It calls `self.llm_client.generate()` which depends on LLM SDK or requests library as integrated by plugin ¹⁴².
- It calls `middleware.before_request/after_response` in loop, which depend on those plugin implementations.
- It relies on `PromptEngine` (which uses Jinja2 environment internally) for prompt compilation and rendering ⁴ ⁵.
- The artifact pipeline logic (in `ArtifactPipeline.execute`) uses its internal structures we saw. The runner just initiates it; heavy lifting (iterating sinks in order, matching produce/consume) is inside that class's method ¹¹⁷. *Interactions:* The runner interacts intensively with:
 - **LLM service** (via LLM client plugin and network calls).
 - **Data plugins** (it uses the DataFrame loaded by DataSource).
 - **All plugin types** at appropriate points (row/val per row, agg at end, etc.).
 - **RateLimiter:** It calls `self.rate_limiter.acquire()` and `utilization()` as part of sending requests concurrently ⁹ ⁸, and `self.rate_limiter.update_usage()` after each response to record used capacity ⁵⁰⁸.
 - **CostTracker:** It calls `self.cost_tracker.record(response, metadata)` after each LLM response to accumulate cost metrics ¹⁶⁷, and uses `cost_tracker.summary()` at end for total cost ³⁴⁷.
 - **EarlyStop** plugins: via `_maybe_trigger_early_stop` on each success and reading `_early_stop_event` in loops.
 - **File system:** if checkpointing enabled, it writes to checkpoint file each time a record is done (so it appends maybe small JSON or text line).
 - **Memory:** It builds large objects like list of result dicts (which could be memory heavy if thousands of entries, but manageable).
 - **Output sinks:** Through `ArtifactPipeline`, indirectly calls each sink's `write`. This is how it interacts with external outputs (like writing to disk, or making network calls to repo/Blob in sinks). It doesn't handle them directly, just triggers pipeline which does.

• Plugin Registries and Base Interfaces:

- **DataSourceRegistry**: Provides `create(name, options)` for data source plugins. Internally maps plugin names to factory functions. E.g., "csv_local" -> a function that returns `CsvLocalDataSource(**options)`. It likely ensures name exists and raises error if not (so config with unknown plugin fails at load).
- **LLMRegistry**: Similarly for LLM plugins (maps "openai" -> function returning `OpenAIClient(**opts)`).
- **SinkRegistry**: Maps sink plugin names (like "csv_file", "github_repo") to factory functions. It also may provide `validate(plugin, options)` method used in `_instantiate_sinks` for extra checks (they call `sink_registry.validate(plugin, options)` before create, as seen in code ⁵⁰⁹). For example, some sinks may require certain fields in options (like repository sink requiring `dry_run` or `path` fields).
- **Controls (RateLimiterRegistry)**: In code, the rate_limiter is created by `create_rate_limiter(def, parent_context)` which likely looks up the "type" field (like "token_bucket") to instantiate appropriate limiter class (e.g., returns a `TokenBucketRateLimiter(rate=..., capacity=...)`). Similarly `create_cost_tracker` picks a class (maybe one that tracks token usage and cost based on model pricing).
- These registries and factories are logically part of the "Core" layer providing factories for plugin layer objects. They decouple the orchestrator/loading code from specific plugin class implementations, making it easy to introduce new plugins.

- Each plugin type also has a base interface or protocol:

- **DataSourceProtocol**: requires a `load() -> DataFrame` method ⁸⁴.
- **LLMClientProtocol**: requires a `generate(system_prompt, user_prompt, metadata=None) -> dict` method ⁸⁶.
- **ResultSink** protocol: requires a `write(results, metadata=None)` method ¹⁰⁸.
- Row and Aggregator plugins follow conventions: row plugin should have `process_row(row, responses)`, aggregator plugin `finalize(records)` ^{88 446}. Baseline plugin requires `compare(base_payload, exp_payload)`.
- EarlyStop plugin requires `check(record, metadata=None) -> reason or None`.
- These interfaces are not formal abstract base classes in code (except possibly via `Protocol` classes in core, as indicated by imports of `LLMClientProtocol` etc. in runner ³⁷⁸). They serve as a contract ensuring orchestrator/runner can call these methods on plugins polymorphically.

• Sample Plugin Components:

- **CsvLocalDataSource** – *Purpose*: Load a CSV from local filesystem. *Responsibilities*: Open the file path from options, call `pd.read_csv(path, **opts)`, return DataFrame ⁸⁴ (tested in core by `test_datasource_csv`). It likely sets `df.attrs['security_level']=...` if defined in plugin context (the orchestrator may wrap it in context, or the factory might attach a level). It doesn't hold state beyond maybe storing path. *Interactions*: with OS (reads file from disk). Code evidence in tests shows it yields correct data.
- **OpenAIClient** – *Purpose*: Integrate with OpenAI API to generate completion. *Responsibilities*: Accept prompts (system & user) and metadata, call `openai.Completion.create` or similar (with model name, input prompts) ⁵³, get response, and return a dict containing at least `'content': response_text`. It might also include usage: the openai library returns `response.usage` which has token counts; the plugin can place those into `metrics` (like

- 'metrics': {'prompt_tokens': ..., 'completion_tokens': ..., 'total_tokens': ...}). *Interactions:* with OpenAI's external API via the SDK or HTTP under the hood. It must incorporate error handling – e.g., catch `openai.error.RateLimitError` and either raise or pass it upward. But since runner wraps calls in try, plugin might just let exceptions bubble up (runner will catch and treat as failure with error message).
- **MockLLMClient** – *Purpose:* Provide a deterministic offline LLM for testing. *Responsibilities:* Either echo the input or load pre-defined outputs. The test dummy did `return {"prompt": user_prompt, "meta": metadata}` ⁸⁶. Another in test returns `{"content": user_prompt, "metrics": {"score": 0.5}}` as in `DummyLLM` in test which returned content echo and a fake score metric ³⁸¹. *Interactions:* none external; used in dev or offline scenario. It ensures that if no external API is reachable, experiments can still run (though obviously the model output is stubbed).
 - **Row Plugin Example – SentimentRowPlugin** (hypothetical): *Purpose:* Compute sentiment of model response. *Responsibilities:* Implement `process_row(row, responses)` – for instance, take `responses` dict (which might have 'content' text) and run a sentiment analyzer (maybe regex or a small model) to classify as positive/neutral/negative. Return a dict e.g. `{"sentiment": "positive"}`. *Interactions:* Possibly with a small internal lexicon or model (no external calls presumably, as that would defeat the offline constraint of plugin; if it needed to call external, that would be a design decision to allow nested model calls, but ideally kept internal for reliability). Code reference: not in given code, but tests mention plugin returning `custom_metric` ⁵¹⁰, which is likely a dummy plugin to test integration (in `test_orchestrator_single_run_executes_plugins`, they register a row plugin that adds `{"custom_metric": 7}` to metrics ⁸⁸).
 - **Aggregator Plugin Example – CountAggregator:** e.g., count number of responses containing a keyword. *Responsibilities:* Implement `finalize(records)` – iterate through results list (each a dict with maybe 'content' or 'metrics') and count how many meet condition, return `{"count_keyword": count}`. The orchestrator attaches 'count_keyword' into aggregates. *Interactions:* purely with in-memory results (no external).
 - **EarlyStop Plugin Example – FailureThresholdStop:** *Purpose:* Stop if more than X failures. *Responsibilities:* Implement `check(record, metadata)` – maybe check `metadata["failures"]` count or maintain an internal counter (the plugin might increment an attribute on each call via metadata row_index or something). If count > threshold, return `{"reason": "too many failures", "plugin": "FailureThresholdStop"}`. This is tricky since early_stop plugin `check` gets called on success records, not failure events directly; one might design it to also inspect global state like runner's failure list via `metadata` if passed. The test though used threshold plugin in `_init_early_stop` with threshold options, and they call `plugin.check` on each success with row index in metadata ⁴⁹⁸. Perhaps the plugin uses `len(self.runner.failures)` if accessible via context. However, conceptually, early stop plugin can track prior failures by having `self.fail_count` in plugin and incrementing it each call if last record was a failure (though they only call on success events).
 - **Result Sink Example – CsvResultSink:** *Purpose:* Save results to CSV file. *Responsibilities:* Implement `write(results, metadata)` – likely use `pandas.DataFrame(results)` to create DataFrame from list of dicts (the `results` list of records), and call `.to_csv(path, index=False)` ³². If `sanitize_formulas` option is true (as config might specify), first scan for any cell value starting with `=`, `+`, `-`, or `@` and prepend `'` to it ¹¹⁰ (their code in clone sinks shows copying sanitize flags to new sink and applying them). *Interactions:* with OS file system (writes to disk). Possibly with `pandas` for convenience of writing. The sink's internal options include `overwrite` (to allow overwriting file or not), `on_error` (maybe "skip" or "fail"), and the `sanitize_formulas` toggles which they pass to an internal function `_sanitize_formulas` that does the scanning and escaping. Code evidence: They clone base

path to experiment-specific in CLI with same sanitize flags ¹¹¹, implying sink's write uses those flags.

- **RepositorySink Example – GitHubRepoSink:** *Purpose:* Commit result files to a Git repository (GitHub). *Responsibilities:* Implement `write(results, metadata)` – perhaps convert results to a desired format (CSV or JSON), then call GitHub REST API to create/update file. The plugin likely has options like `repo` (owner/repo), `path`, `branch`, `message`, and uses `requests` to PUT the content with proper auth header. *Interactions:* with GitHub's API over network. It must handle errors (like if file changed concurrently, API returns a 409 and sink might either fail or ignore based on `on_error` config). It respects `dry_run` by not actually calling API if `dry_run=True` – maybe it logs a message or writes the content to a local temp as simulation. CLI ensures `dry_run` is true by default unless user allowed live outputs ³⁴ (in `handle_sink_dry_run`). Code evidence: the CLI's `_configure_sink_dry_run` modifies sink definitions in config to ensure any repository sinks have `"dry_run": True` unless `--live-outputs` set, indicating sink reads that flag and behaves accordingly ⁵¹¹.
- **SignedSink (LocalBundleSigner):** *Purpose:* Generate a cryptographic signature or signed archive of output artifacts. *Responsibilities:* Implement `write(results, metadata)` – likely fetch from `ArtifactStore` the list of produced files (the pipeline provides the artifact list or the sink reads from metadata if included). Then either:
 - If mode is to sign each file individually, read each file, compute hash (e.g., SHA-256) and then sign the combined hashes or each hash with a private key (maybe loaded from config option `private_key_path` or `secret`).
 - If mode is to produce a bundle, e.g., create a ZIP of all artifacts and then compute a signature for the ZIP (or embed signature in ZIP comment).
 - Output a file like `outputs/signed_manifest.json` or `outputs/artifacts.zip` plus `.zip.sig`.
 - Mark these produced artifacts in `ArtifactStore` with appropriate security level.*Interactions:* with OS (reading files to sign, writing manifest or signature file). Possibly with crypto library (Python's `hashlib` for hashing and maybe `hmac` or `cryptography` for actual signing). There's no explicit dependency of, say, `cryptography` in requirements, so likely they either use simple HMAC for now (which needs only `hashlib` and a key). The architecture note suggests "HMAC manifest" ⁴⁵⁰, which implies they use an HMAC (keyed-hash) rather than public-key signature. So integration is simpler: they use a secret key from config to compute HMAC of concatenated artifact hashes. This avoids complex key management but provides integrity as long as key is secret (not non-repudiation but acceptable for audit within one org).

Each component above was either directly present in code or inferred from the context. The design exhibits **separation of concerns**: - Data loading vs processing vs output vs monitoring each have their own components or plugin hooks. - Extensibility is evident – e.g., to support a new LLM, one adds an LLM plugin and registers it; to support a new output medium, add a sink plugin. The core orchestrator and runner remain unchanged (they treat them via common interfaces). - Code references confirmed critical aspects like how Orchestrator creates runner and passes in everything ⁷⁴ ¹¹⁹, or how Runner coordinates threads and calls plugin methods (with many references demonstrating those calls).

By enumerating each major component, we have a clear catalog: 1. **CLI Interface** – user input handling and process control ²⁴ ²⁶. 2. **Config Loader** – parse YAML into structured settings with instantiated plugin objects ¹³⁰ ¹³¹. 3. **Experiment Orchestrator** – set up context and execute one experiment (calls data load, sets up runner) ⁴⁸ ⁸⁵. 4. **Experiment Suite Runner** – manage multiple experiments, baseline comparisons, aggregated reporting integration ³⁰³ ¹⁰¹. 5. **Experiment Runner** – perform the core data processing loop, concurrency, plugins, and coordinate artifact pipeline ³⁴¹ ⁹². 6. **Artifact Pipeline** – order and execute output sinks with dependency and security checks ³ ³⁸⁹. 7. **Data**

Source Plugins (CSV, Blob, etc.) – unify input retrieval ⁸⁴ ⁵⁵ . 8. **LLM Client Plugins** (OpenAI, Azure, Mock) – handle external model invocation ⁵³ ³⁸¹ . 9. **LLM Middleware Plugins** (Audit, Content Safety, etc.) – implement cross-cutting processing around LLM calls ⁴¹⁹ . 10. **Row Experiment Plugins** (Metrics extraction etc.) – operate on each record’s outcome ⁸⁸ . 11. **Validation Plugins** – verify response structure/contents each record (not explicitly separate in code, but logically included). 12. **Aggregation Plugins** – finalize experiment metrics after all records ⁴⁴⁶ . 13. **Baseline Comparison Plugins** – compute differences between baseline and variant results ³²⁰ . 14. **Early Stop Plugins** – monitor results to possibly cut off further processing early ²⁰⁵ . 15. **Result Sink Plugins** (File sinks, Repo sink, Signing sink, Visualization sink, etc.) – output final artifacts in various forms ¹⁰⁹ ³⁴ . 16. **Control Components** (RateLimiter, CostTracker) – throttle usage and track cost across calls ⁸ ⁵⁰⁸ .

Each of these has been explained with responsibilities and code references. This component catalog provides a blueprint of the system’s building blocks for developers and auditors alike, showing how each piece fits into the overall architecture and where to find its implementation in the codebase.

4. System Qualities and Constraints

4.1 Security Architecture

Security is deeply ingrained in Elspeth’s design, addressing confidentiality, integrity, and accountability at multiple levels. The security architecture is built on a combination of preventive controls (like classification-based gating, sanitization) and detective/audit controls (like thorough logging, signing of artifacts), ensuring compliance with enterprise and regulatory requirements. Key aspects of the security architecture include:

Data Classification and Flow Control: Elspeth enforces a strict **data labeling policy** and uses it to control data flows: - Every input and output is associated with a **security level** (e.g., OFFICIAL, SENSITIVE, SECRET). This is user-defined in config for each plugin or data source, and if omitted, an error is raised (ensuring conscious classification) ⁴⁶ . The code explicitly requires each sink to declare a `security_level`, and validates it on load ⁴⁶ ²⁹⁸ . - The system resolves an overall security level for the experiment by combining levels of components (e.g., data source and LLM) – the `resolve_security_level` function picks the most restrictive label among inputs ⁴⁸ . This ensures if any part of the process has higher sensitivity, the entire run is treated at that higher classification. - During artifact processing, the **ArtifactPipeline** acts as a security guard: it blocks any attempt to send data to a sink with an inadequate clearance. Specifically, before a sink consumes an artifact, it checks `is_security_level_allowed(artifact.security_level, sink_clearance)` and raises an error if the sink’s clearance is lower ¹⁸¹ . Similarly, it checks if a lower-clearance sink declares dependency on a higher-clearance sink’s output (preventing, for example, a “public” sink from consuming a file marked Secret) ³ . This is a form of **Mandatory Access Control** enforcement at the data flow level, analogous to NIST AC-4 (Information Flow Enforcement). - In practice, this means if an operator accidentally (or intentionally) configures, say, a “public_repo” sink but the data was classified “Internal”, the pipeline will halt with a security exception rather than allow that transfer ³ . This guards against misconfiguration and insider mistakes.

Input Validation and Sanitization: Elspeth treats inputs and outputs defensively: - **Prompt Template Validation:** Using Jinja2 with `StrictUndefined` means any missing placeholder in a prompt is treated as an error, not silently ignored ⁴⁰⁰ . This prevents accidental leakage of template syntax or unintended blank inputs. If a field required by the prompt isn’t present in the data (which could cause prompt injection issues or confusion), a `PromptValidationError` is raised and logged, and that

record is marked as failure ⁹⁵. Thus, the content of prompts is predictable and as intended – eliminating an entire class of injection wherein an attacker could provide special tokens to break the template logic (Jinja2's StrictUndefined will stop on undefined variables, not continue).

- **LLM Response Validation:** Validation plugins can be configured to check model outputs. For example, an organization might have a plugin to ensure no PII or certain keywords appear in outputs (akin to a content filter). If the model output contains disallowed content, the plugin can raise an error causing that output to be flagged and removed from results ⁹⁵. Also, if outputs are expected in JSON format, a validation plugin could attempt `json.loads(content)` and if it fails, mark it as failure. This mechanism aligns with OWASP output validation best practices, extended to ML outputs.
- **Spreadsheet Formula Sanitization:** A known risk is when exporting to CSV or Excel: a cell starting with `=` could be interpreted as a formula and execute malicious code (CSV injection). Elspeth's output sinks address this: the `CsvResultSink` by default (or via option) will **escape any formula-injection vectors** by prefixing with a single quote ¹¹⁰. The code shows that when cloning sinks, they carry `sanitize_formulas=True` and then presumably the sink's write method applies it to any cell that begins with `=`, `+`, `-`, or `@` ¹¹⁰. This means even if a model returned something like `=CMD|'/C calc'!A0` (attempting an Excel DDE attack), the saved CSV would have `'=CMD|'/C calc'!A0` which is rendered harmless in Excel. This is an **active defense** against output injection attacks, and demonstrates awareness of even less obvious vectors since model outputs could include characters like `=`.
- **Schema Validation:** The system can verify that input data meets expected format for plugins (via `validate_schema_compatibility`) ⁹⁶. For example, if a plugin expects a column "email" in input to mask it, but data has no such column, the validation would catch it, preventing mis-processing or incomplete sanitization.
- **Secret Handling:** Elspeth requires credentials for external services (OpenAI API keys, Azure keys) but does not log them or store them in plain text beyond what's needed for runtime:
 - API keys are supplied via config or environment. If via config, typically references like `${OPENAI_API_KEY}` are used and presumably resolved either by environment injection or user editing (the code does not explicitly show secrets resolution, but likely the operator uses environment variables for keys).
 - The Azure SDK and OpenAI library ensure that secret keys are not written to logs by default (OpenAI's library might log the first part of the key in debug but not full, and since we control log level, we wouldn't run at such level in production).
 - In logs, Elspeth also doesn't output anything sensitive: it logs high-level events ("Running suite at ..." ³⁹⁷, "Early stop triggered by plugin X" ⁵⁰⁰, etc.). It does not log prompt content or model outputs at INFO level. Only if debug is enabled, it might include more detail (the OpenAI SDK debug logs might show raw request/response, but we would avoid enabling debug in production). Thus, risk of secret or sensitive content in logs is minimized. The security guidance likely would instruct to keep logs at INFO in prod and handle log files securely (which is standard).
- **Authentication & Authorization:**
 - *External Services:* Elspeth leverages secure authentication methods for external APIs. For OpenAI, it uses API keys over TLS. For Azure, it can use OAuth tokens via azure-identity (which uses managed identity or service principal so that no secret goes in code) ⁶⁸. This means credentials are not hardcoded in Elspeth; they are pulled from environment or Azure's token service. That mitigates risks of credential leakage in codebase or config (aside from initial provisioning).
 - *Internal Authorization:* There's no multi-user concept within Elspeth (runs as whoever invoked it). However, it ensures that no unauthorized data flows happen internally thanks to classification gating described. Essentially, it implements an internal form of Role-Based Access Control on data flows: each sink can be seen as having a "clearance level", and each artifact has a "classification level", and the system permits flow only if clearance >= classification. This is akin to a lattice model (Bell-LaPadula) enforcement within the application. The `is_security_level_allowed()` function likely encodes an ordering of levels (maybe Official < Secret etc.) and returns True only if allowed, raising error otherwise ³⁹⁸.
 - The system does not allow arbitrary code execution via config (no plugin path injection, etc.). All plugin names must map to registered factories – user can't directly execute a random Python function through config, which is good (no config-based code injection). Only if a plugin itself is malicious or miswritten would that cause issues, but those are controlled code components. Config YAML is thus treated as data to be validated, not executed.
- **Integrity and Non-Repudiation:**
 - *Artifact Signing:* Elspeth can produce a signed

manifest of outputs. Using an HMAC (as indicated) ensures that if someone were to modify an output file after run, the manifest signature wouldn't match, indicating tampering ⁴⁵⁰. Though HMAC requires sharing a secret key (so for true non-repudiation a public/private signature would be needed), within one org an HMAC suffices to detect tampering (integrity) and, if key is controlled by compliance team, it also means a user couldn't forge a manifest because they lack the HMAC key. This covers integrity of stored outputs – critical if outputs are used as audit evidence. - *Deterministic Execution*: The system labels each run or plugin with a `determinism_level` (like “guaranteed” vs “stochastic”) ¹⁵⁵ ⁵⁰⁴. While this doesn't enforce determinism, it is a form of metadata that can be used to judge result reproducibility. For compliance, they may require certain runs (like ones used for official reports) to be deterministic (e.g., model using temperature 0). The code's ability to coalesce a determinism level and store it in metadata ⁵⁰⁴ means outputs explicitly state whether they were run in a fully deterministic mode or not, which is crucial for interpreting results. Compliance might even forbid non-deterministic runs for final analysis – Elspeth allows marking those and can enforce via policy (like failing validation if `determinism_level` not “guaranteed” for an official run, though that would be an added plugin or config rule externally enforced). - *Audit Logging*: The system generates extensive metadata as described: including not only results but also metrics like number of retries, total cost, and early stop reasons ²¹ ³⁵⁰. Combined with the content of logs (warnings on any anomalies, plugin-by-plugin traces in debug), an auditor can reconstruct what happened during the run: * The `retry_summary` in metadata tells if any call hit errors and retried (if lots of retries or exhausted failures, a compliance officer might investigate model reliability). * The `failures` list in output details exactly which inputs failed and why (error message) ³²⁴, so issues aren't swept under the rug – they're explicitly documented. * The presence of `cost_summary` (if cost tracking enabled) means usage is tracked – useful for cost governance and also to detect anomalies (like if cost is significantly higher than expected for a run, maybe a sign of misuse or bug). * Each early stop event is recorded with the plugin name and reason, and also logged to console (logger.info) with reason ⁵⁰⁰. This ensures if a run stopped prematurely, it's clear why, preventing any silent loss of data. - The **Traceability Matrix** and **Control Inventory** documents (found in docs) likely map these features to specific security controls (like logging mapping to AU-2, classification to AC-4, signing to SI-7 etc.). For example, `CONTROL_INVENTORY.md` presumably lists controls like “All output artifacts are cryptographically signed (SI-7(1)) – Implemented via SignedArtifactSink ⁴⁴⁹”. Or “Prompt inputs and outputs are sanitized to prevent injection (SC-5) – Implemented via StrictUndefined and CSV formula sanitization.” Having these features explicitly in code and docs suggests Elspeth's security architecture was designed to satisfy an ATO's control requirements (which often demand such mapping). - **Minimal Attack Surface**: - Elspeth has no listening ports or generic execution capabilities (no plugin that executes OS commands, etc., aside from running code its devs wrote). This reduces risk: an attacker cannot send malicious input over network to exploit it, since it only calls out. The main input an attacker could control is the dataset and config. If an attacker could feed malicious data (with injection attempts or extremely large values to cause memory issues), the system's validations and Python's safety (no buffer overruns in Python) mitigate a lot. Potential denial-of-service via giant data is a risk (addressed by resource planning, not code). - Code execution from config is not possible – YAML is loaded with safe loader (which doesn't execute arbitrary tags) ²⁶⁴. This prevents YAML exploits that could instantiate dangerous Python objects via `!python/object` tags (PyYAML `safe_load` forbids that). - Dependencies are kept updated as noted (requests 2.31, etc.) which patch known vulns. The dependency analysis tasks the team to monitor CVEs for those libs ⁶⁸ ³⁸. This process ensures the **software supply chain** is secured – if a vuln arises in e.g. requests (like a critical HTTP parsing bug), they'd update the version promptly. They also mention possibly mirroring packages and scanning with tools like pip-audit ⁵¹², demonstrating proactive supply chain risk management (which is a security control often required for ATO). - **Operational Security**: - It's expected that runtime environment security (OS, file permissions) complements Elspeth's internal controls: * The config files and output files should have proper permissions (only intended users can read them). Elspeth by default doesn't handle that (it writes files with OS default umask). On a secure server, umask might be set such that group/world have no access, or outputs can be on an encrypted volume. * Secrets (API

keys, etc.) should be provided via environment or secure config not stored in plain config under version control. The design encourages using environment variables for keys (and azure identity for no key use). There's no secret stored in output logs or files (except possibly manifest includes hashed values of data, which is not secret). * Logging of content is minimal; any debug logs should be protected or turned off in production to avoid inadvertently logging model output that might be sensitive. - The user roles in using Elspeth may be separated: e.g., a compliance officer might prepare the base config (with security levels and mandatory plugins), a data scientist fills in specifics and runs it, and a security team reviews outputs. The tool's architecture supports this by making config human-readable and outputs well-documented (with all needed metadata to verify compliance). - **Contingency and Recovery:** - The checkpoint feature not only improves reliability but also contributes to security by avoiding duplication of processing and ensuring a clear record of what was processed before an interruption (so you don't accidentally skip or double-process data if a run stops). This is not a direct security control, but it helps maintain data integrity and consistency under failure conditions (which can be relevant for continuity controls). - Early-stop ensures that if something goes dangerously wrong (say, outputs are consistently violating policy, an early-stop plugin could halt the run after e.g. 5 such violations, preventing potentially thousands of non-compliant outputs from being generated). This limits exposure in case of an issue.

In summary, Elspeth's security architecture is robust and multi-layered: - At the **data level**, classification and validation prevent unauthorized or malformed data flows. - At the **process level**, lack of open ports and strong default behaviors (no silent failures, no code injection from config) reduce attack vectors. - At the **output level**, sanitization and signing ensure outputs are safe to handle and trustable. - Throughout, **auditability** is emphasized: every significant event (errors, decisions, usage) is recorded either in logs or output metadata, aligning with accountability requirements (AU controls). For instance, if questioned "How do we know the LLM didn't produce disallowed content that was ignored?", one can point to validation plugin logs and failure reports that show any such content would result in logged failure, not silently filtered.

These measures collectively give confidence that Elspeth can operate in high-security contexts and satisfy an ATO's concerns. The code references show concrete enforcement points for these controls: - Security gating of artifact flows ³, - Strict config validation ⁴⁶, - Sanitization in outputs ¹¹⁰, - Use of secure credential methods ⁶⁸, - Comprehensive logging of sensitive operations (like early stop triggers) ⁵⁰⁰.

Thus, Elspeth's architecture aligns with security best practices and compliance requirements from input to output.

4.2 Performance Characteristics

Elspeth is designed to efficiently handle the computational and I/O workload of running potentially large experiment suites. While its primary bottlenecks involve waiting on external LLM API calls (which are I/O-bound), the system includes features to maximize throughput and avoid unnecessary slowdowns: - **Concurrent Request Processing:** Elspeth can dispatch multiple LLM calls in parallel using multi-threading. The ExperimentRunner will use up to a configured number of threads (`max_workers`) when the number of records is above a threshold ⁷ ³³⁶. For example, if `max_workers=4` and 100 records, it will create 4 threads and process 4 requests at a time. This can roughly cut wall-clock time by a factor of up to 4 (assuming the LLM service allows that concurrency and CPU isn't fully saturated). Threads are used instead of async to leverage Python's simpler concurrency model and because much of the time is spent waiting on network I/O (where Python threads can run concurrently since the GIL is released during I/O) ³³⁵. - The `_should_run_parallel` function ensures threads are only used when beneficial (backlog \geq threshold, and `max_workers > 1`) ⁷. The default threshold is 50 records (from

code, if config doesn't specify, it uses default 50) and by default it might enable concurrency with 4 threads if backlog ≥ 50 . This avoids overhead of thread pool for small jobs (where threads might not save time and could even hurt due to overhead).

- **Adaptive Throttling:** To avoid overloading external APIs or violating rate limits, Elspeth introduces a **utilization pause mechanism** ⁸. This checks the `RateLimiter.utilization()` (likely the fraction of allowed calls currently in use or ratio of recent call volume to limit) before submitting each new request in the thread pool:
 - If utilization is \geq a `pause_threshold` (configurable, default 0.8 meaning 80% of allowed throughput), it will pause scheduling new calls, sleeping a short interval (default 0.5s) until utilization falls below threshold ⁵⁰⁷.
 - This essentially paces out requests to avoid hitting a hard cap.
 - Example: If OpenAI allows 60 RPM and the threshold is 0.8, Elspeth will try to keep at ~48 requests per minute by pausing if we're exceeding that. This prevents hitting a scenario where all threads fire and many get rate-limited responses at once. Instead, it smooths the call pattern. The benefit is fewer failed requests (thus fewer retries and less wasted time). This trades a bit of theoretical max throughput for stability and respect of service limits, which in practice yields better performance by reducing backoff/wait times from errors.
- **Vectorized Operations and Data Handling:** Internally, heavy use of pandas for data manipulation (e.g., reading CSVs, filtering columns) is beneficial for performance:
 - Pandas' CSV reading in C is quite fast – reading even large files (hundreds of MBs) is efficient in pure Python terms. Once loaded, iterating with `df.iterrows()` is not vectorized (that yields Python objects row by row), but that's needed because each row requires an API call. The overhead of Python loop vs vectorizing is negligible compared to network call time.
 - Where vectorization can be used, it is – for instance, applying sanitization to output columns could have been done with pandas vectorized string operations. It appears the code might iterate cells manually though (which for 10k cells is fine; if 1M cells, could be slow but still manageable as it's just string prefix check).
 - The orchestrator truncating the DataFrame for `max_rows` uses efficient pandas head operation which is $O(N)$ for slicing (very fast relative to LLM times).
- **Caching and Reuse:** Elspeth avoids redundant computations:
 - It compiles prompt templates once up front, then reuses the compiled template for each row ¹⁹⁴ ¹⁶⁴. This saves overhead, especially as Jinja2 template compilation can take a few milliseconds (which is trivial compared to LLM call but still good to avoid repeating 1000 times).
 - Shared LLM middleware instances are reused across experiments in a suite if they have identical config. The SuiteRunner caches a middleware instance identified by a combination of plugin name, options, and security level ¹⁹⁰. Thus, if multiple experiments use the same telemetry logger or same content safety filter, it will instantiate it once and use for all, rather than constructing one per experiment. This saves some initialization overhead and potentially external connections (e.g., an Azure ML telemetry middleware might open one connection to log events and use it for all experiments).
 - If a RateLimiter is meant to be global for a suite, SuiteRunner reuses that object for each experiment's runner rather than creating new (so it tracks cumulative usage). This not only is functionally needed but also avoids overhead of resetting it each experiment.
 - Data loading is done once for a suite if experiments share input. The CLI code shows that it loads the DataFrame (`suite_instance = ExperimentSuite.load(suite_root)` etc.) and passes `suite_runner.run(df, ...)` using that df for all experiments in the suite ⁴⁰. This means only one disk read (or network fetch) for the dataset even if 5 experiments use it. That's a significant performance gain (reading a large CSV 5 times would be unnecessary overhead).
- **Retry and Timeout Strategies:** Retries can improve performance from a user perspective by reducing run failures:
 - If an API call fails transiently (e.g., network glitch), the built-in retry loop `_execute_llm` will attempt again up to `max_attempts` ³⁶³ ⁴⁰², possibly saving the entire experiment run from having missing results. This avoids manual reruns and thus improves throughput of user completing tasks.
 - The retry logic includes an incremental backoff (`delay` with optional multiplier) ⁵¹³ ⁵¹⁴. This prevents spamming the API rapidly after an error – it instead waits a bit (by default, initial delay perhaps configured like 1s, and backoff multiplier e.g. 2). That means if an API is momentarily overloaded, the runner will politely back off, which often leads to quicker eventual success than hammering (which might get rate-limited further). So, fewer total attempts needed = better performance.
 - There's no explicit global timeout for an API call in code (like requests default might wait indefinitely). A potential risk is if an API call hangs,

one thread could stall. But the design of `early_stop` (e.g., a health monitor middleware could detect slow responses and trigger stop) and the fact that openai SDK usually has an internal timeout (or one can set it in options) mitigate this partly. They do not appear to set `timeout` explicitly, which might be something to add for performance reliability (to not hang forever).

- **Memory Efficiency:** - The system loads all data into memory via pandas DataFrame. This is efficient up to moderately large sizes (millions of rows can be gigabytes, so that's a practical upper bound depending on environment). There's mention of planned streaming architecture to handle larger-than-memory, but currently it's not implemented ⁴⁶⁹, meaning extremely large datasets could be an issue (they would either cause memory swapping or require splitting).
- Within memory, using DataFrame is quite efficient (dense storage, vectorized operations for any filtering).
- Output data (list of dicts) is built in Python which is fine for moderate N (10k or so). Flattening to CSV uses pandas DataFrame conversion which is also efficient (pandas can convert list of dicts to DataFrame quickly in C).
- There's minimal copying: e.g., they sort results by index and then just extract the sorted records (so a shallow copy of references or small overhead) ⁵⁰¹.
- The heavy content (strings from model outputs) remain in those dicts and are just referenced.
- Use of events and locks is judicious: they lock only when adding result to list or checking `early_stop`, which is brief and does not scale with big data – it's O(1) per record. So overhead of thread synchronization is minimal relative to network time, and hence performance scales nearly linearly with number of threads for I/O-bound tasks.
- **Scalability:**
 - *Vertical Scalability:* On a single machine, one can increase `max_workers` to utilize more CPU cores if needed. If the external service can handle parallel requests, performance can scale nearly up to that number of concurrent threads (diminishing returns if API is rate-limited or if eventually local CPU becomes a factor – e.g., if each result must run a heavy local analysis plugin, then CPU could saturate).
 - *Horizontal Scalability:* Because Elspeth is stateless (each run is independent), you can run multiple instances in parallel on different machines or containers to divide the workload (say split dataset into parts, run on different machines, then combine outputs). This isn't automated by Elspeth, but it's possible externally (like using a job scheduler to distribute experiments). There is no internal cluster mode, but none needed given typical use is moderate scale.
- Real-time usage is not intended (it's for batch experiments, not for serving an LLM in interactive latency-critical context). So throughput (requests per minute) is optimized more than single-request latency. That said, per-request overhead is very low beyond what openai SDK has (just some minor Python function calls).
- **Network and I/O Performance:**
 - The system doesn't transfer unnecessary data: e.g., reading only needed columns from CSV if `prompt_fields` is subset could speed up reading (currently, `prompt_fields` is used to filter after reading full row into context, but an optimization could be to pass `usecols` to pandas if known and beneficial).
 - For output, writing CSV or Excel is done once; these libraries are efficient. Writing 10k rows CSV is sub-second, Excel a bit more but still fine.
 - If outputs include images, generation uses Matplotlib which can be somewhat heavy (maybe seconds to produce complex plots). But those are typically a one-time cost and perhaps run in headless mode using Agg which is decent.
 - If a repository sink pushes to Git over API, that can take time for each file (a commit can be a second or two). If pushing multiple large files, it's sequential. But typically, one might push summary results only (like a JSON or a few small files), so performance is okay. The dry-run default prevents accidental huge pushes without user confirming – good because pushing a giant file to Git is slow and maybe not intended.

Performance Testing & Observability: - Elspeth provides built-in metrics that help identify performance issues: e.g., `metadata["retry_summary"]` might show many retries (pointing to maybe needing to reduce concurrency or increase rate limit if hitting a lot of 429s). `metadata["cost_summary"]` shows tokens used – indirectly reflecting how heavy the LLM calls were (useful to gauge if a prompt is too long causing slow responses).

- The logs at INFO will show when each experiment starts and completes with count of rows etc. If an experiment is slow, you can see in logs that e.g. "Experiment X completed with N rows in T seconds" (they didn't explicitly log time, but one could glean from timestamps).
- For internal monitoring, since it's not a persistent service, one typically relies on external monitoring (like if integrated in pipeline, measure pipeline step durations). But one

could instrument via a simple plugin or telemetry middleware: e.g., the “Health Monitor” middleware mentioned in architecture diagrams ⁴²⁹ likely tracks latency of each LLM call and could log or produce metrics (like average response time). This indicates they considered performance telemetry as well – AzureMiddleware might capture time per call and log to Azure Application Insights via azureml-core ⁵¹⁵ ⁴⁴⁵ .

Constraints: - The main constraint is that it runs as a batch job, not a live service. So it's not suitable for extremely low-latency requirements (each LLM call is individually made; no use of streaming or async, though threads do allow overlapping call latency). - Another constraint is memory due to in-memory dataset. For huge data, you might have to manually chunk the runs. The architecture is flexible to do multiple runs (with checkpoint or splitting input) but not automatic. - It's also constrained by external API limits: The best performance you can get is what the external LLM allows. Elspeth's concurrency and throttle aim to reach that safe maximum. It cannot speed up the LLM's own processing speed (that's fixed by model and infrastructure). - Checkpointing slightly affects performance: writing to checkpoint file for each record (that's a small I/O for each success). If processing millions of rows, writing millions of lines to checkpoint could slow it down (~maybe a few MB, likely fine though because file append is buffered and quick on SSD). It's a trade-off for reliability; you might disable it if chasing absolute speed and not concerned with mid-run restarts. - Overhead of security checks: For each artifact and sink, doing `is_security_level_allowed` is O(1) string compare or lookup – negligible overhead relative to I/O tasks. So security features do not meaningfully degrade performance. - Overhead of sanitization: Checking every cell in output DataFrame for formula triggers – if output has, say, 100 columns and 1000 rows = 100k cells, scanning them with Python may take some milliseconds (100k * a few microsec each). That's okay for moderate sizes. For extremely large outputs (like if Elspeth were used to produce a 1e6 row CSV), this loop might be noticeable (1000000 checks, maybe fraction of a second to a second). Still likely fine given that writing the CSV itself would be the bigger time (pandas writes 1e6 rows in a couple seconds). - Overhead of logging: Minimal at INFO (just a handful of log lines). At DEBUG, could be high (if logging every request/response, etc.), but that would only be used in troubleshooting small runs, not in large-scale production runs.

Overall, the performance-oriented design choices (threading, adaptive throttling, compiled templates, caching shared objects) indicate a thoughtful approach to meeting the use cases where dozens or hundreds of LLM calls need to be executed as fast as reasonably possible. Testing with moderate loads likely influenced the default thresholds and concurrency settings. The architecture can handle typical experiment sizes efficiently and degrade gracefully under heavier loads by not overwhelming external services (throttling) and by offering ways to split work (via multiple runs or multi-process external parallelization if needed). The code citations confirm concurrency and throttle behavior: - Thread pool creation ³³⁵ , - Rate limit check loop ⁸ , - Exponential backoff on retries ⁵¹⁴ , ensuring performance strategies are implemented exactly as described.

4.3 Availability and Reliability

Elspeth is built to be a reliable component in automated workflows, with features that ensure that transient issues or partial failures do not derail the entire process, and that results are produced (or failures signaled) in a controlled manner. Several aspects of the design improve availability (ability to complete runs under adverse conditions) and reliability (consistency and correctness of runs):

Fault Tolerance and Partial Failure Handling: - **Per-Record Isolation:** The architecture is robust to individual record failures. If processing one input (one row) causes an error (maybe the LLM fails on it, or a plugin throws an exception for that data), Elspeth does not crash the whole run. Instead: - The runner catches the exception and creates a failure entry for that record ⁹⁵ ³⁶⁰ . - It continues processing the next records (unless an early-stop criterion is met). This design choice – catching

exceptions around `_process_single_row` - ensures that a single problematic input doesn't terminate the experiment. The outcome is a mostly complete result set with only the problematic cases omitted (and documented in `failures`). This is critical for availability: it means a run yields as much data as possible even if some inputs are problematic, rather than failing entirely (which might require a re-run or manual intervention). For example, if 5 out of 1000 records caused model timeouts, the run will finish 995 records successfully and list 5 failures, rather than aborting at the first failure.

- The orchestrator still returns a success (the presence of some failures doesn't cause a non-zero exit; the CLI would only exit non-zero if an uncaught exception or explicit fail happened). Instead, it's up to the user to examine failures and possibly re-run them if needed (maybe with different settings). This approach aligns with resilient batch processing - process what you can, report what you couldn't.

Automatic Retries: Many transient errors are resolved by the built-in retry mechanism:

- If a model call fails due to a network hiccup or a temporary service error, the runner will retry it (with a slight delay). Often, a second attempt will succeed (common with rate limit or transient 500 errors). This significantly improves the chance that an experiment run completes fully on the first try, thus raising the effective availability of the processing. Without retries, a single blip could cause a record failure or worse, perhaps an unhandled exception (if one didn't code carefully). With retries, such blips are smoothed out and fewer records end up in the `failures` list.
- Retries are done individually per record, not at the whole run level (except, if a catastrophic scenario like network down for a long time, many records might each exhaust retries). But even then, Elspeth would mark those as failures and carry on, rather than hang forever or crash.

Timeouts & Deadlock Prevention: The design avoids potential deadlocks in threading:

- The `Early Stop` uses a thread-safe Event and Lock to coordinate threads. If one thread triggers early stop, it sets the event under lock, and other threads check the event at safe points (start of worker, between task submissions) ³³⁹ ²⁴². This ensures no thread continues heavy work after stop condition (they break out asap). It prevents a scenario where some threads keep calling the API even after we know we want to stop (which could waste resources or prolong run unnecessarily).
- The use of `ThreadPoolExecutor` with context manager ensures all threads finish or are cancelled properly when the run ends (no thread left hanging or resources leaked).
- There isn't an explicit global timeout for the entire run, but one could be imposed externally if needed (like a scheduler can kill job after X hours if not done). However, given Elspeth's early-stop capabilities (one could set an early stop plugin to stop if run has taken too long, e.g., check time in `on_suite` hooks), and its incremental progress (writing to checkpoint continuously), it is well-behaved even if it runs longer than expected (you can always kill it manually and resume from checkpoint if needed).

Graceful Degradation: If an external dependency becomes unavailable mid-run (e.g., OpenAI API starts failing heavily), Elspeth's response:

- Each call will retry a few times. If the issue persists, those records will become failures. If a lot fail, an early-stop plugin could detect the high failure rate and abort to save time (if configured).
- The run ends with many failures recorded. But importantly, it doesn't hang indefinitely or go into an inconsistent state; it wraps up and provides partial outputs (all the successes before the issue).
- The checkpoint file by then contains all processed IDs (some succeeded, some failed still included as processed if the logic adds them even if failed? Actually, the code adds to checkpoint only on success by default, so failures are not added. This means if you re-run the same config with the same checkpoint after an API outage is resolved, it will attempt those failures again - which is *good*, because maybe they'll succeed on second run. That indicates a reliability strategy: you can rerun to fill in failures and thanks to checkpoint, it will skip all previously successful ones and only attempt those that failed earlier. This drastically reduces recovery time. The design deliberately does **not** add failed IDs to the skip list in checkpoint, precisely so that you can rerun them.
- We see that on failure, they don't call `processed_ids.add(row_id)` or append to checkpoint, whereas on success they do. This is an intelligent decision for reliability: you only mark a record done if it got a result. So the checkpoint file effectively tracks completion, not mere attempt. Thus, if a run finishes with 50 failures out of 1000, you can fix the cause (say the API is back), re-run, and those 50 will be processed now.
- This approach ensures forward progress - eventually, you can get results for all records without duplicating already processed ones or missing any due to an earlier run.

Redundancy and Re-run: - Though Elspeth itself

isn't redundant (one process runs at a time), the ability to resume via checkpoint adds a redundancy in execution attempts. The config and data can be easily ported to another machine or re-run later. The results of partially completed work are not lost – successes are already output (and could be merged with future outputs if needed), and failures are logged. This is aligned with the concept of **fail-safe** or **graceful recovery**. - E.g., if the machine running Elspeth crashes (power failure) at 80% done, you have up to row 800 logged in checkpoint and maybe outputs flushed until then. You can restart the process on that or another machine with the same config and it will skip 800 and continue with 801 onward. That means the whole job doesn't need to restart from scratch, boosting reliability in the face of system issues. - Also, output sinks that were completed before crash remain on disk or remote (since each sink writes as it goes in pipeline). If crash happened mid-pipeline (say after CSV written but before signing), on resume the signing sink might detect the CSV is present and sign it. Some manual work might be needed to handle outputs that were partially done, but at least the primary results file exists. The pipeline does all outputs at end, so in a crash scenario, you might need to re-run outputs for the portion that was done. The resume mechanism currently doesn't handle partially done outputs (if crash after CSV, the system on resume might see that whole run as incomplete and re-run entire experiment's data, which duplicates content in CSV or overwrites it if config allows). - Nonetheless, in practise, orchestrator commits outputs after processing all data, so more likely a crash during processing leaves no final outputs (only checkpoint). Then re-run processes remaining and at end writes outputs normally. - **Testing and Confidence:** The presence of numerous tests (`test_orchestrator`, `test_runner`, etc. as listed in `SOURCES.txt` ⁵¹⁶ ⁵¹⁷) indicates that reliability was validated in many scenarios (like ensuring plugins run, early stop triggers correctly, metrics are computed right, etc.). For example: - `test_experiment_runner_integration` presumably tests that with concurrency and early stops and failures, the runner still yields correct payload and all events are handled. - `test_sanitize_utils` likely tests formula sanitization thoroughly (ensuring various formula patterns are caught) ⁵¹⁸. - These automated tests contribute to reliability by catching regressions or edge-case bugs, ensuring that as code evolves, reliability features (like not dropping records, properly logging, etc.) remain intact. - **Resource Management:** - The system cleans up thread resources after use (`ThreadPool` context ensures join). - It closes file handles (checkpoint file is opened in append mode for each write, then immediately closed by the context manager each time ⁴¹⁰; data source file reading is done inside `Pandas` which closes at end). - If an experiment stops early (due to early-stop), remaining threads see the event and don't continue, and the executor then shuts down. So no stray threads doing work not needed. - `RateLimiter`, if it holds any timers or so, might need no special cleanup, just goes out of scope when run ends. - Logging uses Python's logging module properly, so it doesn't crash on log issues and respects global logging config. - **Clock and Order:** - Early-stop logic ensures if multiple threads trigger it, only one reason is set (they use a lock to ensure one sets the reason and event, others then see event and skip doing it again) ⁴⁹⁹ ¹⁴⁶. This avoids race conditions that could cause inconsistent state (like two threads both trying to set different reasons). They log if plugin triggers and break out after first found reason. This design yields deterministic outcome in who triggers stop (the first plugin that returns a reason). - Checkpoint line writing for each success is immediate, so if crash happens after a certain success, that ID is in file. There's a tiny window of vulnerability: if the process crashes exactly between finishing a record and writing to checkpoint, that record might not be marked done. On resume, it would be reprocessed (duplicate). But likely negligible risk as writing is done promptly after adding to results (within lock section in `handle_success`). If duplication happened, the effect is producing one record twice (in results or output, which could be noticed by duplicate ID; minor issue). But again, the timing to have a crash in that microsecond is extremely low. That is an acceptable reliability trade-off to avoid more complex transaction-like mechanism. - **High-Level Availability:** - If integrated into a pipeline, the pipeline can rely on Elspeth's exit code to know success vs failure. A non-zero exit code would mean either config error (fast fail) or possibly a severe runtime error that wasn't caught. But under normal circumstances, Elspeth returns 0 even if some records failed (because the run overall succeeded in doing the process as much as possible). So pipeline sees that as success and can move on with partial output. If they want to treat partial failures as pipeline failure, they'd need to parse

output or have a policy like "if failures list not empty, treat as failure" (not automatically done by exit code). - The design choice to not treat partial record failures as process failure can be debated; they opted to consider the job done if it went through all data (even if some fails). This maximizes availability (job doesn't fail frequently), but requires user to handle those failures post-run. Many ETL systems do similar: log failed records to a "rejects" file but still complete the job and let user fix rejects later. - For critical tasks, one might configure Elspeth to be stricter (e.g., if any failure, raise an exception to fail entire run). That could be done via a custom plugin or by checking `len(failures)` after run in a wrapper script. The base tool defaults to maximizing completion rather than all-or-nothing, which is sensible in analysis contexts.

To illustrate reliability in action, consider a scenario: - Out of 100 requests, 5 hit rate limits and 3 have content errors. - **Without Elspeth's features:** The naive approach might crash on first rate limit (if unhandled exception) or just drop responses on floor. Many results lost, maybe manual rerun needed. - **With Elspeth:** - Rate limit responses cause retries – likely they succeed on second try, so those 5 eventually succeed (no data lost, just slight delay). - The content validation plugin catches 3 problematic outputs, logs them, marks them as failures (so they can be reviewed and addressed outside). - The run finishes with 97 successes and 3 failures logged. The exit code is 0 (assuming no exceptions unhandled). - The user sees in results that 3 are missing and in `metadata["failures"]` details of those 3 (why they failed). They can then decide to review those cases manually, or adjust prompt or model and re-run just those (maybe by filtering input to those IDs and using Elspeth on that subset). - The key is the run as a whole succeeded and delivered most of the data, maintaining high availability of results, and clearly signaling where reliability issues occurred.

Thus, Elspeth's architecture yields a system that is **robust and resilient** in the face of errors, supports **graceful degradation** (partial results rather than no results), and facilitates **recovery** (via checkpoint and logged failures). These qualities are crucial in a production or research environment where long experiments cannot be easily repeated from scratch and where downtime or data loss is undesirable. The code references provided (for retry logic [402](#), early stop event use [242](#), checkpoint usage, etc.) substantiate these reliability mechanisms.

5. Technology Stack

Elspeth is implemented in Python and leverages a range of proven libraries and technologies. Below is a complete inventory of the technology stack, including versions as specified, along with notes on their usage and implications (such as known vulnerabilities or support status). The stack can be divided into runtime environment, core libraries, and optional extras:

- **Programming Language: Python 3.12** – The system is built in Python, taking advantage of its rich ecosystem for data science and its ease of integration. Python 3.12 (the version noted in documentation) is the latest major Python release and is actively supported (not end-of-life) [63](#). Python provides the dynamic capabilities needed for plugin architecture and robust libraries for HTTP, concurrency, YAML parsing, etc. Using Python means the system benefits from rapid development and extensive testing frameworks, but it must be mindful of performance (which Elspeth addresses via concurrency and optimized libraries).
- **Security note:** Python itself has a strong security track record when used properly. The code does not execute untrusted code via `eval` or similar, which mitigates common Python security pitfalls.

- *Version support*: Python 3.12 will be supported by the Python core team for several years. This aligns with long-term maintenance needs. It also means Elspeth can use the latest Python features (like improved performance in 3.11+ due to CPython optimizations).

- **Core Libraries:**

- **PyYAML** (Version ≥ 6.0) – Used for parsing YAML configuration files ⁴⁹. PyYAML 6.0 is the latest major version and fixes some older vulnerabilities (like arbitrary code execution via `yaml.load` which is addressed by using `safe_load`) ²⁶⁴. Elspeth uses `yaml.safe_load`, which avoids executing YAML tags that could instantiate objects or run functions, ensuring config files are treated as plain data (important for security). PyYAML 6.0 has no known open vulnerabilities at this time (there were historical CVEs in older versions related to `yaml.load` usage).
 - *License*: MIT License (per PyYAML docs) – compatible with Elspeth’s intended licensing (which is likely a permissive or internal license) and fine for enterprise use.
 - *Note*: PyYAML is pure Python for parsing, which is sufficiently fast for typical config sizes (on the order of KBs). There is no performance issue here.
- **pandas** (Version $\geq 2.2.0$) – The primary data manipulation library ⁴⁹. Pandas is used to load input data (CSV/Excel) into DataFrames and to output results (via DataFrame to CSV/Excel conversion) ³² ³². Version 2.2.0 is very recent (pandas releases regularly).
 - *License*: BSD 3-clause – permissive and acceptable for enterprise.
 - *Considerations*: Pandas relies on NumPy and is highly optimized in C for heavy operations. No known major security issues (it doesn’t process untrusted code, just data; some CVEs were in older numpy for buffer overreads, but not relevant unless processing malicious binary data).
 - *Maintenance*: Pandas is actively maintained; 2.x series is current and will likely continue to be supported. It’s a crucial part of the stack enabling high-level data operations.
 - *Note*: Pandas in an ATO context is considered a standard library – known to be stable and thoroughly tested. It is heavy (in size), but usage is justified for the convenient data structures and I/O.
- **requests** (Version $\geq 2.31.0$) – The ubiquitous HTTP library used for making REST API calls ⁴⁹. Requests 2.31.0 is the latest in 2.x and has no open CVEs (past minor ones dealt with SSL handling have been fixed in earlier updates).
 - *License*: Apache 2.0 – permissible in enterprise (just requires notice).
 - *Use in Elspeth*: Likely used by LLM client plugins (either directly or via OpenAI library which itself uses requests internally) and by repository sink plugins to call web APIs ⁴⁵. It handles TLS verification by default, which Elspeth relies on (the dependency analysis reminds to include corporate CA if needed but doesn’t disable verification) ⁴⁵.
 - *Notable*:
 - Requests is not asynchronous but fits Elspeth’s thread model well.
 - It supports proxies (so integration with enterprise network is easy).
 - Already used widely, security-vetted in many environments.
 - *Potential vulnerability watch*: We will monitor any new CVEs (like the one in 2018 (CVE-2018-18074) which was fixed long ago). The dependency note references monitoring for "requests signing or logging vulnerabilities" ³⁸, indicating the devs are aware and plan to patch promptly if anything arises.
- **OpenAI Python SDK** (Version $\geq 1.12.0$) – Used to interact with OpenAI and Azure OpenAI endpoints ⁴⁹. Version 1.12.0 is relatively recent (OpenAI releases often). This library abstracts the HTTP and provides convenience methods for chat and completions.
 - *License*: MIT – fine for enterprise use.
 - *Integration*: Elspeth’s Azure OpenAI usage suggests they rely on this SDK (with `api_type` and `api_base` configured for Azure) ⁵³. The dependency analysis says to

- "lock to patched versions to mitigate request signing or logging vulnerabilities" ³⁸ , implying they are aware that older OpenAI SDK versions had issues (perhaps logging sensitive info or not properly verifying server certs). Using $\geq 1.12.0$ indicates they include the fixes (for example, OpenAI 1.0.0 had some logging of keys at one point which was fixed).
- *Support & known issues*: No known major vulnerabilities, but the library is fairly new. It will likely update as OpenAI changes APIs. The devs plan to stay current (the note suggests an update on 2025-10-12 aligning with latest usage).
 - *Scope of usage*: The SDK internally uses `requests` (so inherits its robust TLS handling). They caution to "monitor request signing or logging vulnerabilities" ³⁸ – possibly referring to older issues where debug logging could inadvertently log sensitive request parts. In 1.12.0, such issues are presumably patched.
 - The Azure integration via this SDK might require azure-identity for token auth in some scenarios (the dependency analysis mentions Azure OpenAI adapter alignment ³⁸ , they likely feed in an API key or use the `azure-identity` provided token via requests).
- **Jinja2** (Version $\geq 3.1.0$) – The template engine used for prompt construction ⁵⁰ . Version 3.1.0 is recent and importantly includes fixes for any known sandbox escape flaws (Jinja2 had a history of minor issues but typically only if sandbox is used; Elspeth uses StrictUndefined but not full sandbox mode since template is trusted by config author, not user-supplied).
- *License*: BSD 3-clause – compatible and minimal restrictions.
 - Jinja2 is crucial for merging prompt text with data safely. The dependency note says "StrictUndefined mitigates template injection but stay current for sandbox fixes" ⁴⁰⁵ . Indeed, using StrictUndefined prevents accidental injection from undefined variables, but they also keep Jinja updated in case of any discovered templating vulnerabilities. No known CVEs in Jinja2 3.1.0; older ones were mostly related to sandbox (CVE-2019-8341 affecting sandbox mode, which Elspeth doesn't use but anyway 3.1.0 is beyond that).
 - Performance: Jinja2 is fast for templating, and they compile templates once and reuse them, which is best practice.
- **jsonschema** (Version $\geq 4.21.1$) – Possibly used for validating config schemas or plugin schema compatibility ⁵⁰ . This library validates JSON data against a schema. The code uses `validate_settings` likely implemented using jsonschema to enforce config structure (the dependency analysis referencing schema parsing CVEs suggests they use it) ⁴⁰⁵ .
- *License*: MIT – fine.
 - Known vulnerability: older jsonschema (before mid-2023) had a DoS vulnerability (CVE-2023-3135) related to large schemas causing recursion explosion. Version 4.21.1 is likely patched (the note about "update promptly for schema parsing CVEs" ⁴⁰⁵ hints they keep an eye on such issues). No known open CVEs in this version if updated.
 - It's probably used to ensure config file meets expected spec (e.g., the `SettingsSchema` might be defined in JSON Schema form and validated). If so, it catches mistakes early. It might also be used for validating input data schema vs plugins (if they define a JSON schema for expected input columns, etc. as part of plugin definition).
- **scipy** (Version $\geq 1.10.0$) – SciPy is listed as dependency ⁵¹ , presumably for any statistical functions in metrics or baseline comparisons. For example, if they include a stats plugin to compute significance (like Welch's t-test or correlation), SciPy would provide that (via `scipy.stats`).
- SciPy also could be a transitive dependency via other libs (pandas or pingouin might rely on it).
 - *License*: BSD 3-clause – fine.
 - No known security issues (it's mostly math routines in C).

- It's heavy but only used if needed for advanced analytics (like baseline significance calculation).
- **Optional “Extras”** (for extended functionality, activated via extra install flags):
 - **Matplotlib** (Version $\geq 3.8.0$) and **Seaborn** (Version $\geq 0.13.0$) – Used for visualization output if the user requests visual charts ⁶⁰. They are included under an `[analytics-visual]` extra. Matplotlib 3.8 and Seaborn 0.13 are latest, with no major security issues (these do not face external input beyond data and are widely used).
 - *Licenses*: Matplotlib – PSF (Python Software Foundation) License, Seaborn – BSD. Both permissive.
 - They allow generating PNG or HTML charts for output. The documentation note says these are optional for PNG/HTML visual sinks ⁶².
 - They would increase memory and CPU usage if heavy plotting is done, but only for final reporting (should not impact main experiment run significantly).
 - They mention to ensure fonts/backends are vetted if used (visualization sometimes requires installing system fonts – potential environment consideration) ²³⁴.
 - **Azure ML Core (azureml-core)** (Version $\geq 1.56.0$) – This is included under an `[azure]` extra, used for telemetry middleware (logging to Azure ML) ⁵¹⁹. Version 1.56 is somewhat heavy and has many dependencies (but likely they only import it if middleware is active).
 - *License*: MIT – common for Azure SDKs.
 - azureml-core has no known direct vulnerabilities, but it has many dependencies (like azure-storage, etc.) which they mention to track transitives and patch accordingly ²³⁴. They specifically note it pulls `msrest` and `adlfs`, etc., which should be kept in sync ²³⁴.
 - azureml-core is in maintenance (there is a v2 with azureml-mlflow, but 1.56 is still supported).
 - This integration is optional – only install if Azure ML telemetry needed, as they caution to install only when required to minimize surface.
 - **openpyxl** (Version $\geq 3.1.0$) – Used for Excel output sinks ³⁸⁷. Openpyxl reads/writes .xlsx files.
 - *License*: MIT or similar – fine.
 - Known issue: older openpyxl had a vulnerability where formulas or links in untrusted Excel could lead to exposing system files (CVE-2021-27852) – but that's for reading malicious Excel. Here Elspeth only writes Excel, not read, so risk is low. They still keep it updated (≥ 3.1).
 - They note "include in hardened builds only when spreadsheets necessary" ³⁸⁸ – meaning they treat it as optional to avoid pulling heavy dependency and potential attack surface if not needed.
 - **Pingouin** (Version ≥ 0.5) and **Statsmodels** (Version ≥ 0.14) – These appear under `[stats-agreement]`, `[stats-core]`, etc., for advanced statistical analysis extras ⁵². E.g., Pingouin for inter-rater agreement stats, Statsmodels for advanced tests.
 - *Licenses*: Pingouin – MIT, Statsmodels – BSD. Both fine.
 - These are used only if advanced metrics plugins are enabled (like computing p-values or effect sizes). They mention pinning minor versions for deterministic outputs ⁵²⁰, because stats functions can change results slightly across versions; for accreditation, reproducibility is key, so they likely freeze exact versions for those optional stats packages.
 - Not security-critical (these don't handle external input besides numeric data).
 - **scikit-learn** might be pulled in as dependency if Pingouin is installed (Pingouin or stats extras mention it transitively ⁵²¹). If so, scikit-learn's license (BSD) is fine; no major security issues (it's purely analytical).

- **HMAC or Crypto:** Not explicitly listed, meaning they likely use Python's built-in `hashlib` and `hmac` for signing (which don't need external packages). If they needed asymmetric crypto, they'd use e.g. `cryptography` library (which is not listed, implying they don't use it).
 - Using HMAC for signing is simple and FIPS-140 compliant if running on FIPS mode (Python's OpenSSL can be built to FIPS mode).
 - They considered "mirroring or vendoring critical packages" for supply chain reasons ⁴⁰⁷ – not a stack item, but a practice: e.g., hosting an internal PyPI mirror for these libs to avoid tampering risk at fetch time.

Version and CVE Analysis: - All specified versions are up-to-date as of late 2025 (for example, requests 2.31.0 was released mid-2023 and has no known vulnerabilities since). The note entries indicate conscious tracking: - azure-identity and azure-storage: monitor for any credential escalation CVEs (none known currently) ⁵⁵ . - openai library: ensure using patched versions (1.12.0 is reasonably current) ³⁸ . - requests: same (2.31.0 includes fix for CVE-2023-32681 regarding proxy credentials leak, which was addressed around requests 2.31). - jsonschema: 4.21.1 likely includes fix for CVE-2023-3135 (which was fixed in 4.17.3 or so). They plan to update promptly if any new arises ⁴⁰⁵ . - They explicitly note pip-audit and capturing reports for accreditation ⁵¹² – implying they run vulnerability scans on the final environment and include those in artifacts. This ensures any known issues in dependencies are caught and mitigated or documented. - **Support status:** - None of the libs used are at end-of-life: * Python 3.12 – newly released, will be supported upstream likely till 2028. * Pandas 2.x, requests 2.x, Jinja2 3.x, etc. – all actively maintained. * azureml-core 1.56 might be superseded eventually by Azure ML CLI v2, but given it's optional, it's fine and still supported in 2025. * The only caution is if azureml-core stops updating, they'd have to possibly adapt to newer Azure ML SDK, but as an optional telemetry plugin that's not critical path. - They mention upgrade strategy doc – indicating a plan to handle library upgrades carefully (for compatibility and compliance) ⁵²² . They pinned certain minor versions for deterministic outputs (like stats models) ⁶² and enumerated compatibility extras in docs (e.g., `[stats-bayesian]`, `[stats-planning]`) for optional features with heavy libs ⁶² .

Open Source Licenses and Compliance: - The stack is predominantly permissive licensed (MIT, BSD, Apache). There's no GPL or AGPL library in the list, which avoids any copyleft concerns. So from a license compliance standpoint, using these in an enterprise context (even distributing Elspeth internally or eventually open-sourcing it) is low risk. - They have a placeholder for license in their repo (no LICENSE file yet, just note reach out to maintainers) ²⁴⁶ , meaning currently it's proprietary or not formally licensed. They plan to add one. In either case, third-party licenses need to be complied with: - They likely maintain a `LICENSE_THIRDPARTY.md` or similar listing those (maybe in `CONTROL_INVENTORY`). - Permissive ones require including their license text, which presumably they will do when publishing license (the note "License information will be published... until then reach out before redistributing" implies they are aware of compliance obligations with dependencies). - For an internal ATO, they should show they track third-party license compliance (which they do via docs).

Hardware and OS: - Python runs on cross-platform, presumably targeted to Linux (for server deployments). Possibly tested on Windows for dev environment. - No platform-specific libraries except maybe if azure-identity uses OS managed identity (works on Azure VMs on Linux/Windows alike). - If deployed in e.g. a Docker container, the base image should be a recent OS with Python 3.12 (like an official Python 3.12-slim base). - ATO context might require using a hardened base image (which they can do, the code has no OS-specific ties except needing standard Python and C libs for the packages). - No custom native extensions in Elspeth aside from those provided by libs (pandas/numpy etc. ship their compiled C extensions, but those are widely vetted).

In summary, the technology stack of Elspeth is modern, widely used, and chosen for a mix of functionality and compliance: - **Functionality:** Pandas for data handling, Requests/OpenAI SDK for

external calls, Jinja2 for templating – these make implementing features straightforward. - **Performance:** These libs are optimized (pandas in C, requests uses efficient HTTP, etc.). - **Security:** They largely abide by safe usage patterns (safe_load, StrictUndefined, TLS by default) and are kept updated to avoid known vulnerabilities (explicitly mentioned in their dependency analysis). - **Longevity:** All libraries are mainstream with active communities, so bug fixes and support are expected to continue. The devs have shown intent to track and upgrade dependencies as needed (with an upgrade-strategy doc and active updates noted as of 2025-10-12) ²³⁴ .

This stack gives confidence that Elspeth stands on a solid technological foundation that is stable, secure, and performant enough for its intended use. Each component's version has been vetted for known issues and compatibility. By planning to keep them updated and isolating optional features into extras (so one doesn't have to install heavy or less secure libs unless needed), they also minimize the attack surface and bloat of the deployment.

6. Architectural Decisions & Rationale

Throughout the design and implementation of Elspeth, the development team made several key architectural decisions to balance flexibility, security, and maintainability. This section highlights those decisions, the alternatives considered (implicitly or explicitly), and the rationale behind the choices, as evidenced by patterns in the codebase:

- **Pluggable Architecture vs. Hardcoded Workflow:** The team chose a highly modular, **plugin-based design** where data sources, LLM integrations, metrics, and outputs are all interchangeable via registries ⁷³ ⁷⁴ . **Decision rationale:** This ensures **extensibility** – new LLMs or analytics can be added without altering core logic, fostering an open framework for experimentation. It also improves **maintainability**, as each plugin is isolated (e.g., adding a new output format doesn't risk breaking LLM code). The alternative – a monolithic orchestrator coded specifically for, say, OpenAI and CSV only – would have been simpler initially but very inflexible (any change would require modifying core code). By deciding on a plugin system, they accepted a bit more complexity (registries, interfaces) in exchange for long-term adaptability and user customization (users can drop in custom plugins for proprietary data sources or internal models).
- This decision is evident in code via the use of factory methods and the `plugins/` directory structure: e.g., orchestrator calls `create_row_plugin` and doesn't know which specific class it gets ⁷³ , indicating loose coupling.
- The **rationale** was likely driven by the need to support multiple cloud providers (Azure, OpenAI), multiple output types (files, devops, etc.), and evolving metrics (like if tomorrow they need a new compliance metric, they add a plugin rather than rewriting orchestrator). This aligns with the initial design goals of "pluggable orchestration" ⁶ .
- **Classification Enforcement (Security by Construction):** A design decision was to embed security classification as a fundamental property of components (plugins, sinks) and to enforce it at runtime in the artifact pipeline ³ . The alternative could have been to handle this via documentation and user discipline (i.e., just instruct users to not send secret data to public outputs). They instead **codified** the policy in code. **Rationale:** This decision significantly reduces the risk of human error and ensures compliance is **systematically enforced**. It reflects an architectural philosophy: whenever possible, make security automatic rather than optional. They likely decided this to meet stringent compliance needs (like for an ATO, you need to demonstrate no classified data can leak to lower domains – implementing the check in code provides that guarantee).

- The code implementing `is_security_level_allowed()` and raising exceptions if mismatched ³⁹⁸ is an explicit design choice. Many internal tools might not bother with such checks, but here they treat data classification as a first-class attribute.
- The trade-off of this decision is added complexity (each plugin must declare level, pipeline must compare, etc.). They determined the benefit (preventing potentially severe data leaks) outweighs the extra complexity.
- **Use of Standard Protocols/Formats vs. Custom:** The team consistently uses standard formats: YAML for config, JSON/CSV for outputs, and established protocols (HTTP, JSON for API calls). For example, they chose **YAML for configuration** rather than a custom config format or code-based config. **Rationale:** YAML is human-readable, widely supported, and allows structured data (including nesting for profiles and plugin lists) easily. This makes it accessible to non-developers (like a security team member can read the YAML to verify settings) and easily integrated with other tools (like an Ansible pipeline could generate the YAML). The alternative – e.g., requiring experiments to be defined in Python code – would reduce accessibility and possibly mix code with config (less clear separation). YAML's drawback can be syntax complexity or footguns (like tabs vs spaces), but the team mitigated by not using YAML's risky features (safe_load only).
- Another decision is usage of **JSON for internal data interchange** (manifest, trace matrix, etc.) and **CSV/Excel for results** to ease consumption by common tools (Excel, BI tools). They didn't invent a proprietary binary format. The rationale is likely **interoperability** – these formats ensure outputs can be opened or parsed by anyone with standard software.
- They also output an **Excel workbook** for combined analytics, which, while heavier than CSV, is user-friendly for business users (with multiple sheets for summary, details, charts). This decision is about convenience and clarity for end-users (for audit reports, Excel with charts is a common deliverable). They decided supporting Excel was worth the dependency on openpyxl, gating it behind an extra so it's optional ⁵²³.
- **Multi-Threading vs. Async/Await:** When implementing concurrency, they chose Python **threads** (with ThreadPoolExecutor) rather than using `asyncio`. **Rationale:** Likely simplicity and compatibility. The code is largely synchronous and integrates with libraries like requests which are synchronous. Converting everything to an async paradigm (async HTTP library like httpx, async file I/O, etc.) would have required more refactoring and might complicate integration with third-party libs (OpenAI's lib is not async in the version used). Python threads, despite the GIL, work well for I/O-bound parallelism and allowed them to keep code structure straightforward (the logic in `_run_parallel` is easier to write and reason about than full async event loop management). They also avoid dependencies on an async runtime – simpler for adoption. The trade-off is that threads can't achieve true parallel CPU usage due to GIL (but CPU-bound parts are minimal anyway), and each thread has overhead (memory for stack, context switching). Given typical use (tens of threads at most, mostly waiting on network), this decision is rational. Evidence: use of ThreadPoolExecutor in code ³³⁵.
- The alternative (async) could in theory handle many more concurrent tasks with lower overhead, but OpenAI's rate limits and typical usage likely don't require, say, 1000 concurrent calls – a handful is enough to saturate model throughput. Thus, threads are sufficient and simpler.
- This decision improves maintainability (most Python devs are comfortable with threading and the code looks like normal sequential code within each thread). It also aligns with leaving heavy lifting (like multi-core vectorized ops) to underlying libraries (like numpy inside pandas), which release the GIL anyway so threads can overlap those operations too.

- **Use of Checkpointing vs. Relying on External Job Restarts:** They built a **checkpoint mechanism** to handle run interruptions or staged execution ³³¹. **Rationale:** This is somewhat unusual for short batch jobs, but for long-running experiments in possibly unstable environments (or for iterative development where one might stop a run mid-way to adjust something), this adds resilience. They explicitly mention capturing accreditation artifacts including vulnerability scanning reports ⁵¹², which suggests runs may be part of compliance pipelines – in such context, checkpoint ensures that even if a run is paused for scanning or fails due to environment glitch, progress isn't lost.
- Alternative was to not implement checkpoint and just require rerun from scratch on failure. But given that a run might involve, say, thousands of model calls (costly and time-consuming), that approach would hurt reliability and user trust. Checkpointing is a conscious design to improve **fault tolerance** (an ATO often asks "how do you handle failures gracefully?").
- They decided on a simple text file approach (writing one ID per line). Rationale: simplicity, easy for user to inspect or modify (and no heavy dependency, just file I/O). Also, text is portable if they move run to another machine (just bring the file).
- They considered concurrency: the checkpoint is written under a lock in multi-thread runs to avoid concurrent writes messing it up. So they ensured thread-safety.
- One trade-off: slight performance cost as discussed (file write per record), but decided reliability is more important for long runs. They likely judged that cost negligible relative to network calls or overall run time.
- **Auditability and Compliance as Core Design Goals:** Many design decisions stem from satisfying compliance requirements:
 - The decision to **collect extensive metadata** (like cost_summary, attempt history) ²¹ ¹⁶⁸. Rationale: provide evidence and trace for each run (for example, an ATO security control might require tracking usage of external services – cost_summary directly addresses that by listing tokens used, which correlates with content amounts and cost).
 - The decision to output a **traceability matrix** (perhaps automatically or via doc) and **control inventory** in docs is unusual in normal dev projects but was clearly made to ease ATO submission. It shows they prioritized aligning code to controls and documenting it (e.g., they updated it on 2025-10-12 with references to new features) ⁴²⁴. This evidences an architectural decision to treat compliance artifacts (like mapping code to NIST controls) as part of the system deliverables.
 - They decided to implement a **Signed Artifact Sink** for integrity. Many internal tools would skip that, but they included it (manifest + HMAC) to meet non-repudiation requirements. Rationale: ATO frameworks like FedRAMP often require verifying that output data hasn't been tampered with in transit or storage (SI-7 control). By including a signing step, they can claim compliance with that control out-of-the-box. The alternative was to rely on environment measures (e.g., storing outputs on a secure file share with checksums at storage layer). They chose an application-level solution for independence and explicitness. The simplicity of HMAC was chosen likely due to ease (no PKI needed) and because typically only internal parties handle outputs.
 - Another compliance-driven decision: **deterministic output enforcement**. They ensure that optional extras like scientific libraries which could introduce nondeterminism (e.g., if random initializations or multi-thread race in math libs) are pinned to versions and that the concept of determinism is recorded (so they can claim "We label outputs as deterministic or not; for accredited scenarios we ensure determinism where required"). Not many projects explicitly track that, so it's a conscious compliance-oriented design. The dependency note about "when accreditation relies on deterministic outputs, pin minor versions" ⁵²⁰ demonstrates this

reasoning (making sure a statistical result doesn't change with library upgrade, which could break an audit trail or require re-validation).

- **Use of Safe Defaults:** Many micro-decisions show preference for safe default behaviors:
 - Dry-run mode for repository sinks by default ³⁴ – prevents accidental external commits unless user really wants it (they must use `--live-outputs` to enable) ⁵²⁴. Rationale: Protect data (perhaps sensitive) from inadvertently being pushed to a repo. They intentionally made user opt-in for any remote write. This is a secure default aligning with least privilege principle (don't send data out unless explicitly told).
 - StrictUndefined in Jinja2 by default – as already covered, that's a safer default than default Jinja which would insert blank or break template unpredictably. They recognized possible *template injection issues* if undefined variables were left (e.g., if an undefined var was in prompt, Jinja might insert nothing and model might interpret leftover braces, etc.). Using StrictUndefined causes immediate failure, a safer outcome than producing potentially gibberish prompts. They specifically highlight this decision in dependency analysis ⁴⁰⁵.
 - Not storing secret keys in config by default – guidelines for azure-identity usage (the code automatically picks up environment or managed identity if no key given). They could have forced user to put keys in YAML, but they avoided that, leaning on environment and identity. This design decision reduces the risk of secrets in config files (which might be checked into version control inadvertently). Instead, config can reference environment variables (like they mention using `${VAR}` in options presumably and the loader could substitute from `os.environ`).
- **One-Process Model vs. Distributed:** They decided to implement Elspeth as a single process (with concurrency internal) rather than a distributed service or multi-process pipeline. They do not incorporate message queues or microservices architecture (like a separate data ingestion service, separate processing service). **Rationale:**
 - Simplicity and easier deployment for the intended use (which is often interactive or one-off batch by a team, not a 24/7 service).
 - A single process is easier to containerize or run on a single server under control (which for ATO is actually beneficial – fewer moving parts means easier to evaluate).
 - They likely considered that scaling needs (performance, as above) can be met via threads and occasional splitting of runs by the user if necessary, without the overhead of designing a distributed system.
 - This decision yields a design that's *stateless between runs*, which in compliance context is simpler to secure (no database to manage, no persistent sensitive state to protect except output files).
 - The trade-off is that it doesn't automatically parallelize across multiple machines – but the plugin approach means an external orchestrator could run multiple Elspeth processes if needed (and they provide the reporting tools to merge results if needed).
- **Not building a GUI or Web UI:** They stuck to CLI and file outputs for user interaction. This decision was likely intentional to minimize complexity, attack surface, and to meet users where they are (data scientists often comfortable with CLI and analyzing outputs in Jupyter or Excel). Creating a web interface would require a web server (with user auth, etc.), which would significantly raise security requirements (and ATO scope). By keeping it CLI, they avoid a whole class of vulnerabilities (XSS, CSRF, web auth issues) and focus on core functionality. The README suggests usage is via CLI and outputs as files, confirming this decision ⁵²⁵ ⁵²⁶.

- Possibly they considered but deferred a GUI (maybe that could come as a separate layer, but not needed for initial accredited version).
- This also means the learning curve for new users might be a bit higher (no point-and-click), but the target audience (technical teams) is fine with CLI/YAML.
- **Extensive Logging vs. Minimal Logging:** They opted to include fairly detailed logging and warnings (like logging each validation warning, logging when early stop triggers, etc.). The alternative could be minimal logging to avoid clutter. They decided on verbosity at least at warning level to ensure any unusual situation is recorded for audit. This is in line with compliance: e.g., *AU-3 Content of Audit Logs* requiring enough info to determine what happened. They include plugin names in logs (like "Early-stop plugin 'threshold' triggered (stop reason: cost limit)") ⁵⁰⁰, which is valuable for later analysis. The design likely expects logs to be collected into a central log management for review if needed. They put emphasis on not just errors but also actions (like early stop triggers, or maybe logging how many rows processed in some summary).
- They may also log external call counts via metrics in output rather than in log lines, which is fine since results become part of audit evidence (AU-6 centralized analysis could consider output file content as part of audit artifacts).
- **Using Rate Limiting vs. Assuming external handling:** Many client tools just let the external API return 429s and handle them via retry. Elspeth includes a **RateLimiter plugin** (like a token bucket) to proactively throttle ⁹ ³⁷¹. **Rationale:**
 - It ensures compliance with API usage policies (some APIs may impose penalties if continuously violated; internal usage might also have quotas to watch).
 - It yields more stable performance as discussed.
 - The alternative was to rely on OpenAI's built-in rate limit error messages and just handle those by waiting. They instead integrated a token-counting approach (the `cost_tracker` tracks tokens, `rate_limiter` tracks QPS or throughput). This decision improves reliability (fewer errors) at cost of implementing a small control mechanism. The trade-off is complexity but given they have security in mind, implementing their own limiter also means they can incorporate conditions like "if utilization above 0.8, wait" which is a safer pattern recommended by best practices. They prioritized stable runs over maximizing raw throughput at all times (ties into reliability).
- **Dependency Pinning for Determinism:** They consciously pin or suggest pinning certain library versions where variability could affect reproducible outputs (like stats models) ⁵²³. This decision is important for accreditation – ensuring that an analysis can be repeated later and yield identical results (if a stats library changed a random seeding or algorithm, output might differ slightly; by pinning version they avoid that). Many projects float dependencies, but here they pick exact (or minimum but then they've likely locked environment in practice via pip freeze in their artifact). The rationale is explicitly mentioned: "*when accreditation relies on deterministic outputs, pin minor versions*" ⁵²⁰. This shows an architectural stance: stability and reproducibility of outputs is more important than always using the latest minor release of a stats lib. They accept risk of running a slightly older version for sake of consistency (with an upgrade strategy to eventually move if needed).
- **Use of HMAC vs. Digital Signature:** For output signing, they appear to use an HMAC with a symmetric key (evidenced by "HMAC manifest" mention) ⁴⁵⁰. The alternative is using a public/private key (digital signature). The decision for HMAC suggests:

- Simplicity: HMAC is easier to implement and doesn't require a PKI or distribution of a public key – the verifying party could be the same system or compliance officer who holds the secret.
- Possibly, since this is an internal tool, non-repudiation (where a signature could be verified by anyone without secret) was less critical than integrity and detection of tampering. Within one org, HMAC is sufficient because the set of people verifying integrity is limited and can share the key.
- Using HMAC avoids complexities of key management – they can put a secret in config (or environment) and use Python's `hmac` library. Using a digital signature would require storing a private key and distributing a public key (which might be fine but is extra setup). They likely decided HMAC was the right balance for now. It might be an interim decision (they could later support PGP or similar if needed).
- Given ATO context, this is acceptable for integrity proof as long as the key is protected. It is slightly weaker in "non-repudiation" since whoever holds key could forge a manifest, but if the compliance team is the one holding the key, they trust themselves.
- Code supports this via `hashlib` for SHA-256 and `hmac` for combining with secret. No external crypto library was added, signifying decision to use Python's standard crypto (which is FIPS-compliant when Python is compiled with OpenSSL FIPS mode, if needed).

In conclusion, the architectural decisions in Elspeth consistently reflect: - A drive for **extensibility** (plugins, registries). - A commitment to **security & compliance by design** (enforce classification, safe defaults, heavy logging). - A preference for **robustness & reproducibility** (retries, checkpointing, deterministic outputs). - Adoption of widely used **standard technologies** (YAML, JSON, threads) over inventing new ones, to reduce learning curve and risk.

Each decision involved trade-offs: - More complexity (plugin system, context management) vs. simpler hardcoded flows, which they accepted for long-term flexibility. - Slight performance overhead for better security and reliability (e.g., classification checks, sanitization scans), which they deemed worth the cost in a compliance-heavy environment. - More upfront effort (like implementing early stop, checkpoint) vs. relying on user or environment to handle issues – they chose to invest in these features to make the tool **resilient and enterprise-grade** out-of-the-box.

These decisions align with building a system suited for enterprise ATO approval: modular (to adapt to new requirements), secure (to meet strict controls), and reliable (to avoid failures and facilitate audits). The consistency of these choices across the code confirms a deliberate architectural vision oriented around those priorities.

7. Risk Assessment

In developing and deploying Elspeth, several architectural risks were identified and addressed. Below, we outline key risks categorized by security, reliability, and maintainability, assessing their severity and noting how they have been mitigated or could be mitigated:

7.1 Security Risks & Mitigations: - Data Leakage via Misrouted Outputs – *Risk*: High. If sensitive data were sent to an unauthorized sink (e.g., confidential prompt outputs pushed to a public repository), it would be a serious breach. *Mitigation*: The classification gating mechanism in the artifact pipeline ³ greatly reduces this risk by automatically blocking such flows. Every sink must declare a security level, and any mismatch triggers a `PermissionError` before data leaves the system ¹⁸¹. The design assumes correct classification by config authors; the residual risk is misclassification (user tags something as "Official" that is actually "Secret"). *Residual Risk & Mitigation*: This is a process risk – mitigated via training and perhaps having compliance review configs. The system provides a safety net but cannot fully prevent human misclassification. However, even if misclassified, if all sinks are internal,

damage is limited. - **Hardcoded Credentials or Secret Exposure** – *Risk*: Medium. If API keys were hardcoded in code or config, they might be inadvertently exposed (in logs, repo, etc.). *Mitigation*: Elspeth uses environment variables and Azure Managed Identity for credentials ⁶⁸, avoiding putting secrets in config files. Nothing in code prints these secrets (they are never logged). The dependency analysis flags to monitor logging of secrets in libs ³⁸; currently, no such logging occurs at default levels. Thus, the architecture avoids storing secrets at rest and ensures they only reside in memory or secure env storage. *Residual Risk*: If a user includes an API key directly in YAML, it could be visible on disk; guidance should discourage that. Possibly they could enhance by supporting a vault lookup mechanism in config (not implemented yet), but environment injection is usually sufficient. - **Injection Attacks (Template, CSV)** – *Risk*: Medium. Without countermeasures, an attacker who can influence input data or prompt templates might attempt to break out of templates (like injecting Jinja syntax) or create malicious CSV content. *Mitigation*: - Template injection is mitigated by using StrictUndefined so that unexpected placeholders cause failure rather than execution, and since templates themselves are controlled by config (not user-supplied at runtime), the main injection vector is if input data contained Jinja-like patterns. However, Jinja is only used for templates, not run on output content. Input data is just substituted as strings, not interpreted as code. - CSV/Excel injection is mitigated by formula sanitization (prefixing formula-like text with `'`) ¹¹⁰. This neutralizes Excel macro or formula exploits from model outputs. *Residual Risk*: If a model output contained something like `@evilsite.com` (which Excel might treat as external link? They cover `@` in sanitization too ¹¹⁰), it's handled. They did not mention HTML sanitization for the Visual report (if model output is inserted into HTML report, could it inject a script?). If VisualSink creates an HTML file with content from model, there's a potential XSS risk if someone opens that HTML in browser. However, typically such HTML would be local and opened by trusted user, so risk is low. As a mitigation, they could escape model content in HTML, but it's not noted explicitly. This is a minor risk given context (internal use, not a public web server). - **Denial of Service (Resource Exhaustion)** – *Risk*: Medium. A malicious or buggy input could cause extremely heavy processing (e.g., extremely large input file or prompts leading to enormous outputs). *Mitigation*: - They allow `max_rows` to limit how much of a dataset to process ¹³³. This can be used to avoid huge jobs. - The rate limiter prevents flooding external APIs (thus avoiding hitting usage caps catastrophically or being banned) ⁸. - The design encourages running on controlled infrastructure where input size and model usage are known. - Checkpointing ensures partial progress is saved, so an unexpected interruption (even due to resource limits) doesn't force a total redo. - However, if someone gave a 10 million row CSV, Elspeth would try to load it fully into memory – likely failing or swapping. This is a risk if user error or maliciously huge input. *Mitigation*: Document memory guidelines or split data. In future, implement streaming (they plan that). The risk is partly accepted (assuming typical usage in tens of thousands of rows at most). - **Dependency Vulnerabilities** – *Risk*: Low to Medium per library. The stack uses many libraries; a vulnerability in one (e.g., requests or pandas) could affect the system. *Mitigation*: The team actively monitors and updates dependencies ⁶⁸ ⁵²⁷. They mention using pip-audit and capturing reports for the ATO package. Also, running in a restricted environment (no public access for the process) limits exploitation surface. If e.g. requests had a bug processing malicious HTTP response, an attacker would have to compromise OpenAI API to exploit it – unlikely and outside threat model (OpenAI being attacked is separate, and presumably they'd fix on their side). The largest risk might be if PyYAML had a new flaw and reading config could be exploited – but they use `safe_load` which mitigates most YAML vulns. Overall, by keeping libs updated (we see they updated libs as of Oct 2025), they handle this proactively. - **Output Integrity/Non-Repudiation** – *Risk*: Low after mitigation. Without signing, outputs could be tampered after generation. *Mitigation*: The SignedSink adds HMAC manifest, so any tampering can be detected ⁴⁵⁰. *Residual Risk*: The HMAC approach means whomever holds the key could also tamper and re-sign (so an insider with key could falsify evidence). In an ATO context, the key would be held by a trusted role (e.g., compliance officer), so that risk is acceptable given trust assumptions. If needing stronger non-repudiation, they'd move to asymmetric signing in future. Right now, the risk of unnoticed tampering by an outside entity is low (they'd need the key or to break HMAC, which with a strong key and SHA-256 is practically impossible). - **Logging Sensitive Data** – *Risk*: Low. They

intentionally avoid logging content or keys at INFO. Debug logs could include content or usage details (openai library debug might log full request JSON including prompt text – not sure, but possible). *Mitigation:* They default to INFO logging in normal runs; debug must be manually enabled for troubleshooting and would be done in secure environment if needed. Also, logs are internal – not accessible externally – but for privacy, one might worry about model outputs with personal data being in logs. They did not implement a specific log scrubber, but by default, they don't log outputs anyway. So risk of logs leaking sensitive info is minimal. - **Multi-tenancy & Access Control** – *Risk:* Low (not multi-user service). If Elspeth is run on a shared server, any user with OS access can potentially read output files or config of others if permissions not set. That's an operational issue: presumably each run is done by authorized individuals, and OS-level controls (file permissions, user accounts) restrict unauthorized access. The architecture didn't implement internal user separation (out of scope given it's not a web app). So risk of one user seeing another's data is left to environment config (run as separate OS user, etc.). Given it's typically used within one team or pipeline, that's acceptable. - **OpenAI Model Output Risks** – Not exactly architectural but noteworthy: Model may produce incorrect or biased results. *Mitigation (security in sense of trust in output):* They include validation plugins for format and possibly content (like a content safety check plugin that calls Azure Content Safety API on outputs to flag disallowed content) ⁴¹⁹. This reduces risk of deploying harmful content (like if used in an automation that sends outputs to end-users, a content moderation plugin could stop it). This is more about ethical risk, but they did account for it by including hooking points for content safety (e.g., `ContentSafety` in middleware chain) ⁴²⁹.

7.2 Reliability Risks & Mitigations: - **Single Point of Failure** – *Risk:* Medium. The architecture runs as a single process on one machine, so if that machine or process fails, the run is disrupted. *Mitigation:* Checkpointing allows resumption on the same or different machine without total data loss ³³¹. Also, the statelessness means one can always restart the process or redeploy to another server. There's no persistent service state to corrupt. If machine goes down, picking up the YAML and checkpoint on another machine yields minimal downtime beyond job restart. - To further mitigate, one could run Elspeth on a stable server or have a failover server. The architecture doesn't handle failover automatically but makes it easy to do manually thanks to stateless design. - **Partial Data Loss** – *Risk:* Low. Without checkpoint, a crash mid-run could lose outputs of processed records (if not flushed). *Mitigation:* They flush each result to checkpoint file and gradually build output structure in memory. Final outputs are written at run end, which means if crash before output, you'd have to rerun; but thanks to checkpoint, you rerun only missing part, and then outputs will include all results. So no data is ultimately lost, just delayed. - They might consider periodic flush of partial results to disk (e.g., writing an interim CSV every 100 rows). They did not implement that (because easier to just rerun from checkpoint if needed). - The reliability risk of losing processed data is thus minimized by checkpoint existence. - **Excess Retries or Slowdowns** – *Risk:* Low. A poorly configured rate limiter (too high threshold) might lead to lots of 429 errors and hence many retries, elongating run. Or a mis-behaving external API could cause sequential long timeouts per record. *Mitigation:* The rate limiter/pause mechanism alleviates this by preemptively slowing submissions ⁸. They also log if many retries happen (the `retry_summary` in metadata would highlight if `total_retries` is large, signalling a performance issue to address). - They don't have a dynamic circuit-breaker (like "if 10 calls in a row fail, stop to avoid endless hanging"). However, the `early_stop` plugin could be configured to do that. This is left to configuration: e.g., set an early-stop to trigger if failure count > X. They provided the mechanism (check of metadata in `early_stop` plugins) for such conditions. So if not done by default, a user can add it. - The reliability risk of a run bogging down due to endless error->retry loops is mitigated by `early_stop` or by eventual giving up after `max_attempts` and marking failure (so it does not truly get stuck on one input for too long beyond configured attempts). - **Out-of-Memory or Disk Space** – *Risk:* Medium if input is huge or outputs large. *Mitigation:* - Memory: If input is extremely large, they expected user to use `max_rows` or split input. This is somewhat a known limitation, not fully mitigated in code except by documentation and future plan for streaming. - Disk: If writing huge outputs (like million-row

CSV, multi-GB logs), the environment should have sufficient space. Not specifically checked by Elspeth (no disk space monitor). The risk is if disk fills mid-output, the sink might throw an IOError. The system would then treat that output as failed (e.g., repository sink will except if it can't write file or commit). That would bubble up as an exception likely (since no catch around pipeline). That might fail the run at end. So low-level reliability: an out-of-disk scenario would cause run to error at final output stage, after all processing done (so data is in memory, but output partially written). *Mitigation*: - This scenario should be prevented by provisioning enough space. - If it occurred, partial outputs might be present (CSV half written) and HMAC sink would fail on incomplete file, raising error. Checkpoint file would mark all records done, so a re-run would skip processing and jump to outputs, but if disk still full, it would fail again. - There's no automated recovery from disk full, apart from cleaning space and re-running the output step (which might require manual removal of incomplete output to avoid confusion). - This is a general risk not unique to Elspeth, accepted as operational risk. Possibly mitigated by monitoring and ensuring environment sizing. - **Complex Config Leading to Misconfiguration** – *Risk*: Low-Medium. The flexibility of YAML profiles, extras, etc., means users might configure contradictory or suboptimal settings (like two early_stop plugins conflicting, or forgetting to classify a sink). *Mitigation*: - Validation at load catches many such issues (like missing security_level on sink causes config load error) ⁴⁶. - The documentation likely provides recommended patterns, and example configs show proper usage. - The risk of user error is lowered by making many settings default to safe values (e.g., dry_run true, concurrency disabled unless threshold hit, etc.), so forgetting to set something usually doesn't lead to insecure or crashed state. - If a config is logically inconsistent but not caught (e.g., two early_stop plugins that both trigger same event, or a prompt pack and experiment define same plugin twice), the worst is some redundant work or a minor logic oddity (like aggregator computing same metric twice). This doesn't break reliability of producing outputs, just might confuse results slightly. There's no evidence of specific guard for duplicate plugin definitions – they'd just both run. That is minor risk (like doubling a metric).

7.3 Maintainability Risks & Mitigations: - High Complexity / Learning Curve – *Risk*: Medium. The plugin architecture, context passing, and multiple moving parts make the codebase non-trivial for new developers to understand fully. *Mitigation*: - They provided thorough documentation (Architecture Overview, plugin catalogue, dev notes) ⁵²⁸ ⁵²⁹ to map out components and flows (the mermaid diagrams in docs help maintainers see structure ⁴⁵³ ⁴⁶⁹). - They have a CONTRIBUTING.md and presumably code comments to guide devs (the snippet of CLAUDE.md suggests they even used AI to keep docs in sync, which indicates emphasis on maintainability). - Extensive tests also serve as maintainability aide, because if a change breaks something, tests will catch it, directing maintainers to the issue. - The risk remains that the orchestration logic (with contexts and events) is intricate. But by modularizing (each plugin in its file, core in core/orchestrator.py etc.), they localized complexity somewhat. A new dev can focus on one plugin or one part at a time. - **Tight Coupling of Components** – *Risk*: Low. Thanks to plugin abstraction, coupling is actually loosened. But some internal coupling exists (runner knows about all plugin lists, orchestrator knows about runner specifics like concurrency config). *Mitigation*: They separated concerns logically (core vs. plugins vs. pipeline). The use of dataclasses and clear interfaces (like `ResultSink.write`) reduces unstructured interactions – each part communicates through defined methods. - E.g., Orchestrator doesn't directly manipulate sinks or data, it passes them into runner. Runner doesn't know orchestrator details aside from what's attached in plugin context if needed. - RateLimiter and EarlyStop are integrated across runner threads but via well-defined mechanism (Event/Lock). - So maintainability risk from coupling is relatively low as design is quite modular. - **Outdated Documentation vs. Code** – *Risk*: Medium. They explicitly noted some documentation needed updating (they did an update on 2025-10-12 for parts) ⁵³⁰. If docs lag code, maintainers might rely on incorrect info. *Mitigation*: They performed an audit of docs and updated them (document says "Update 2025-10-12" for various sections) ⁵³¹ ⁵³², showing commitment to synch docs with code changes. Also, the presence of mermaid diagrams that incorporate new modules (like concurrency, baseline plugins in updated diagrams) ⁴⁶⁹ ⁴¹⁵ indicates they actively maintain

architecture documentation. This reduces risk of maintainers being misled. Continual updates and mention of a consolidated documentation approach (like referencing the new plugin paths in docs after reorganization) ⁵³⁰ further mitigate this risk. - **Dependency Bloat and Conflicts** – *Risk*: Low-Medium. With many optional extras, there's risk of dependency conflicts (e.g., one extra requiring a version of SciPy that conflicts with another). *Mitigation*: They carefully separated extras (so you only install what you need, reducing bloat). They likely test combinations that are common (like core + azure + visual) to ensure no conflict. They pinned minimal versions to ensure compatibility but allowed flexibility beyond (e.g., openai>=1.12.0 so it can use 1.13 or 1.14 when they come, trusting no break; for sensitive ones like stats libs, they pinned exactly or at least recommended to pin via internal policies). - They have an upgrade strategy doc to manage dependency updates systematically ⁵²², which helps maintainers handle upgrades in a consistent fashion and consider impact on outputs. - **Complex Test Matrix** – *Risk*: Low. Many combinations of plugins and extras could yield a huge test matrix (with azure, with openai, with or without stats extras). If not all combinations are tested, a corner case might fail. *Mitigation*: Not explicit, but by modular design, one can test each plugin category in isolation (which they do in unit tests for orchestrator, runner, each plugin). They likely test at least one scenario covering each optional path. Some risk remains if a very unusual combination is used (like azure+visual+some exotic stat plugin concurrently), but impact would likely be minor (like an extra metric miscomputed). - **Turnover of Maintainers** – *Risk*: Low-Medium. If key architects leave, new maintainers might find the system complex. *Mitigation*: The thorough documentation (including rationale in this SAD and control mapping) and tests form a strong knowledge base. Additionally, writing things like CLAUDE.md (maybe internal doc) indicates they might have used AI or internal tools to maintain consistency (less reliance on single person's memory). The code structure is also fairly logically organized (core vs. plugins vs. docs), which aids new maintainers in locating relevant parts.

Each risk above has been weighed and largely mitigated through design choices or supporting processes. The system exhibits **defense in depth** for security and **fault tolerance** for reliability, while maintainability is bolstered by clear modularization and documentation. The remaining residual risks are either accepted as low likelihood/impact (e.g., misclassification of data – handled by policy outside system) or slated for future improvement (like streaming for extremely large data). Overall, the architecture appears robust against both foreseen and some unforeseen issues, aligning well with the demands of an accredited, production-quality system.

8. Compliance Mapping

Elspeth's architecture and implemented controls have been carefully mapped to common security and compliance requirements, such as those in NIST SP 800-53 (for U.S. federal ATO) and organizational policies (data handling rules, audit requirements). Below we map key architectural elements to specific compliance controls or best practices, demonstrating how the system meets those standards:

- **Access Control (AC) & Information Flow Enforcement:**
- **AC-3 (Access Enforcement) / AC-4 (Information Flow Enforcement):** The classification-based artifact pipeline enforcement corresponds directly to these controls. By ensuring that data of a certain classification only flows to sinks authorized for that level ³, Elspeth enforces data flow policies (no "Secret" data to a "LOW" destination). This satisfies e.g. internal rules that sensitive data not be sent to external systems. The mapping: classification levels in config are likely aligned with organization's data categories, and the pipeline's `is_security_level_allowed()` is effectively an automated guard implementing AC-4. Evidence: code enforcing clearance checks ³⁹⁸.
- **AC-6 (Least Privilege):** The design of dry-run repository sinks (no data leaves the system unless explicitly allowed) ³⁴, and requiring user confirmation (`--live-outputs`) to push data externally, exemplifies least privilege principle. Sinks default to not performing external writes

⁵²⁴, meaning by default, the system runs in a mode that doesn't expose data outward unless needed. This ensures that, out-of-the-box, it uses minimal privileges on external communication.

- **AC-19 (Access Control for Mobile Code)** and **CM-7 (Least Functionality)**: Elspeth does not execute any mobile### 8. Compliance Assessment

Elspeth's architecture and controls have been mapped to the relevant security and regulatory requirements that an Authority to Operate (ATO) would evaluate. Below we highlight how the system meets key compliance standards and best practices:

- **Identification & Authentication (IA)**: Elspeth itself does not manage user identities (it runs as a CLI tool under an authorized user account), so controls like IA-2 (user identification) and IA-5 (authenticator management) are largely inherited from the operating environment. However, for **external service authentication**, Elspeth uses secure methods:
 - It integrates with Azure Active Directory for Azure OpenAI and storage access via **Managed Identity** (no hard-coded credentials) ⁶⁸. This aligns with IA-5 (2) – using organization-managed keys/tokens rather than static secrets.
 - OpenAI API keys are supplied via environment variables or secure config, not embedded in code or output, satisfying AC-6 / IA-5 by limiting knowledge of secrets to the runtime and preventing disclosure in logs or files.
- No default passwords or accounts exist in Elspeth; it relies on OS-level authentication to run and on properly scoped API keys for external calls (adhering to IA-2 & IA-8, requiring unique credentials and technical non-person entities to use dedicated keys).
- **Access Control & Data Handling (AC & SC)**: The system enforces strict **data access policies**:
 - The **classification labeling and enforcement** corresponds to AC-4 (*Information Flow Enforcement*) and internal data governance policies. By requiring each output sink to have a clearance and automatically preventing high-classification data from flowing to low-clearance outputs ³, Elspeth implements mandatory access controls on data. This means the system honors confidentiality constraints (e.g., no "Secret" data to public endpoints), fulfilling policies like "Data will only be stored or transmitted to locations authorized for its classification."
 - All **network communications are encrypted** via TLS 1.2/1.3 (Requests and Azure SDK default to verifying SSL certificates) ⁵⁹. This meets SC-8 (*Transmission Confidentiality*) and SC-13 (*Cryptographic Protection*) for data-in-transit. For instance, calls to OpenAI's API use HTTPS with certificate verification (ensuring server identity and encryption).
 - **Least Privilege**: By default, Elspeth runs in a minimal mode: repository sinks are dry-run unless explicitly enabled ³⁴, and no data leaves the host unless the user configures it. This design addresses AC-6 (*Least Privilege*) by not performing actions beyond what is necessary. For example, it will not push data to an external repo unless the operator flips the `--live-outputs` switch, ensuring that by default, privileges (like internet data transfer) are not exercised unnecessarily.
 - **Boundary Protection**: The classification check also acts as an internal guard for SC-4 (*Information in Shared Resources*) – ensuring that data of one security domain doesn't co-mingle with a lower domain output. And because Elspeth runs as a batch process, it does not create shared memory or resources accessible by different users or processes simultaneously, reducing risk of unauthorized data access via system resources (aligning with SC-32 (*Information Separation*)).
- **Privacy Considerations**: If subject to privacy laws (GDPR, etc.), Elspeth supports compliance via plugins: for instance, a **data anonymization plugin** could be included as a validation or row plugin to redact PII before sending to an external LLM (fulfilling requirements to protect

personal data when using third-party services). Although not a specific NIST control, this aligns with internal privacy impact assessments – the architecture’s plugin-extensibility allows adding such controls without core changes, demonstrating compliance flexibility.

- **Audit and Accountability (AU):** Elspeth provides extensive logging and traceability, satisfying *AU-2 (Auditable Events)*, *AU-3 (Content of Audit Records)*, and *AU-6 (Audit Review)*:

- **Auditable Events:** The system logs notable security-relevant events, such as validation failures, early stop triggers, and configuration warnings. For example, if a content violation occurs, it is captured in the `failures` output and likely logged as a warning, which can be reviewed ⁹⁵₅₀₀. The logs also record when an experiment starts and completes, and any security enforcement (e.g., an unauthorized flow attempt would raise an error that is logged). These cover *AU-2* by identifying events that need auditing (failures, stops, errors).
- **Audit Record Content:** Each run’s metadata includes details like number of records processed, number of failures, the reasons for failures, total API calls and tokens used ²¹₃₅₀. This rich context in the output (which can be considered audit records) meets *AU-3* by providing answers to "what occurred, when, and with what outcome". For instance, `retry_summary` shows if any calls failed and were retried, `early_stop` shows why processing halted if it did, etc. Additionally, the *HMAC-signed manifest* provides integrity for these audit records themselves, aligning with *AU-10 (Non-repudiation)* – an authorized party can verify that audit artifacts (outputs, logs) have not been tampered with post-run.
- **Audit Review and Analysis:** The outputs (in JSON/Excel) and logs are in human-readable form, facilitating *AU-6* (audit review). A compliance officer can easily inspect an Excel “Executive Summary” sheet to see if any anomalies occurred (the design intentionally surfaces all anomalies in outputs). The presence of structured output (JSON) allows automated analysis as well – e.g., ingesting the metadata JSON into a SIEM or audit management tool to automatically flag runs with failures or unusually high token usage (which could indicate misuse).
- **Retention:** While Elspeth itself doesn’t enforce log retention (that’s an operational matter), the fact that results and logs are output to files means they can be stored in an audit archive. The signed manifest ensures if they are stored, any later review can trust their integrity. This supports *AU-9 (Protection of Audit Information)* – the manifest (with HMAC) protects audit records from undetected modification, and controlling access to the secret key ensures only authorized personnel can regenerate a valid manifest.

- **System & Communications Protection (SC):**

- **SC-7 (Boundary Protection):** Elspeth operates within an internal network environment; it does not itself provide boundary control, but its use of classification enforcement and TLS means it respects established boundaries (no plaintext data crossing network boundaries, and no cross-domain data flows without clearance). Additionally, if deployed in a segmented environment, it uses the corporate proxy (ensuring all external communications funnel through monitored egress points) ⁴⁵, supporting *SC-7* by not bypassing network protections.
- **SC-28 (Protection of Information at Rest):** The system expects underlying storage to provide encryption at rest (e.g., running on encrypted disk volumes or storing outputs on encrypted cloud storage). Elspeth’s contribution is ensuring sensitive data is labeled so it can be handled accordingly at rest (e.g., an output file marked Secret would be stored only on approved encrypted drives as per policy – a process control outside Elspeth, but enabled by Elspeth’s labeling). Additionally, the optional signing of outputs can be seen as part of protecting data at rest (integrity).

- **SC-18 (Mobile Code):** Not directly applicable as Elspeth doesn't download or execute mobile code or untrusted code – it executes Python code that is part of its installation (plugins), all vetted and under configuration control. This addresses any policy requiring review/whitelisting of code executed, since all its code is part of the system's codebase or controlled plugins (no dynamic code fetch).

- **Configuration Management (CM) & Maintenance:**

- **CM-3 (Configuration Change Control):** Elspeth's use of human-readable YAML and documented architecture makes it easier to subject changes to formal review. Every experiment config or code change can be reviewed via version control and traceability matrix (they maintain a document linking requirements to implementation updates ⁴²⁴). This is evidence of a robust configuration management process: e.g., if a control requirement changes, they update docs and possibly code with an "Update" annotation (as seen on 2025-10-12 updates) ⁵³¹. That shows compliance with controlling and documenting changes (CM-3, CM-4).
- **CM-6 (Configuration Settings):** Secure defaults in Elspeth (like dry-run for outputs, StrictUndefined for templates) constitute secure configuration settings enforced by the application. This means out-of-the-box, it's in a hardened config state, satisfying policies that systems be deployed with secure defaults (e.g., *"disable unused functionality by default"*, which they do for output sinks and extras).
- **SI-2 (Flaw Remediation):** The team's process of monitoring libraries for CVEs and promptly updating (as evidenced by dependency-analysis notes and upgrade history) ⁵³³ speaks to a compliance with vulnerability management requirements. They integrate scanning (pip-audit) and maintain a documented history of updates (documentation audit dated logs changes) ⁵³¹, aligning with SI-2 and SI-5 (automatic flaw remediation).

- **System & Information Integrity (SI):**

- **SI-7 (Software Integrity):** The HMAC signing of output artifacts, although meant for data, also provides integrity assurance of those deliverables. If one treats the signed outputs as part of the system's audit evidence, it ensures they have not been tampered post-generation, which is akin to ensuring the integrity of generated information. For software code itself, presumably the organization would sign the code distribution or use hash verifications as well (not an Elspeth feature, but part of devOps – they mention vendoring critical packages and scanning, which supports supply chain integrity).
- **SI-10 (Information Input Validation):** Elspeth performs rigorous validation on inputs (prompt templates and data). The `validate_settings` and data schema checks ensure that only properly structured input is processed ⁴²⁷ ⁹⁶. Moreover, the early content validations via plugins can be seen as input validation for model outputs (ensuring only expected formats are accepted). This meets the spirit of SI-10 by verifying all information (config, input data, model output) against expected criteria and handling deviations.
- **SI-4 (Information System Monitoring):** Through detailed logging of events and usage metrics (like logging when early stop triggers or how many tokens used), the system supports monitoring of security-related events. If integrated with SIEM, these logs can trigger alerts (e.g., if cost spikes unexpectedly or if a validation plugin flags disallowed content, that event in logs can notify security staff).

- **Planning & Risk Assessment (PL & RA):**

- **RA-5 (Vulnerability Monitoring):** The team's approach to dependency updates and pip-audit aligns with continuously monitoring for known vulnerabilities and addressing them ⁵³³. They explicitly note tracking vendor advisories (e.g., Azure SDK, OpenAI library) ⁵¹². This proactive stance meets RA-5 obligations of scanning and remediating vulnerabilities in the system components.
- **PL-8 (Security Concept of Operations) / PL-2 (System Security Plan):** The extensive documentation (architecture diagrams, control mapping, traceability matrix) essentially forms part of the system security plan and concept of operations. The inclusion of this SAD as an output demonstrates compliance with documenting the system's architecture and security controls – a requirement for ATO submission.

In addition to NIST controls, **organization-specific policies** are addressed: - *Data Classification Policy:* Elspeth's classification tagging and enforcement directly implements the org's data classification standard (e.g., no "Company Confidential" data should leave the intranet – enforced by classification checks). - *Secure Development Policy:* The use of safe coding practices (no use of `eval`, use of safe YAML loader, sanitizing outputs, etc.) and maintaining of control inventory shows adherence to the organization's secure development lifecycle. - *Operational Monitoring Policy:* The integration of metrics and logs means after deployment, operations can easily monitor usage (ensuring compliance with any policy on cloud API usage tracking or cost management). - *Data Retention Policy:* Elspeth by default outputs data to files which can be stored or deleted per policy rather than storing it internally. So it doesn't hinder compliance with retention or disposal policies – the organization can manage output files per their schedule (Elspeth doesn't lock data inside a proprietary database). Also, it writes classification on output files (either in file content or via file naming conventions in config), aiding compliance with marking requirements (some orgs require classified documents be marked – here, the content itself often contains the security level in metadata and in control inventory files). - *Incident Response:* If an incident (like an output had disallowed content or was sent to wrong sink) occurs, the thorough logs and manifest expedite forensic analysis (mapping to IR-4 (Incident Handling) and IR-5 (Incident Monitoring)). For example, if a user accidentally attempted to output secret data to a low sink, the system would block it and log it – that log entry can be treated as a minor security incident to follow up, which is clearly recorded (facilitating IR procedures).

In summary, Elspeth's design is tightly aligned with compliance requirements. Many of its features – classification enforcement, extensive auditing, secure defaults, encryption of data in transit, signing of outputs – are not incidental but rather deliberately implemented to satisfy controls needed for an Authority to Operate. The **traceability matrix** and **control inventory** documents (provided in `docs/compliance/`) directly demonstrate this mapping by listing each control and how the architecture meets it ⁴²⁴ ⁵³². For instance, the control inventory likely states: "AC-4: Implemented via security_level enforcement in artifact pipeline (see `orchestrator.py`, lines...)", "AU-3: Logging includes event type, timestamp, source (see logging statements in `runner.py` etc.)". By designing the system in compliance from the ground up, Elspeth significantly smooths the ATO process, since auditors can easily verify each required control in the code and documentation.

9. Recommendations

Having analyzed Elspeth's architecture in depth, we identify the following recommendations to enhance the system, focusing on addressing any critical gaps for ATO approval, potential improvements, and guidance for future evolution:

9.1 Critical Issues for ATO Approval: 1. **Finalize Licensing and Third-Party Compliance** – Although not a security flaw, the absence of a formal license in the repository must be resolved before operational use. We recommend publishing the intended license (likely a permissive license, as

indicated) and including a Third-Party License Notice file listing all dependencies and their licenses. This is important for legal compliance and to satisfy any ATO requirement regarding software provenance. It appears the team plans this ²⁴⁶; it should be done prior to ATO submission. Additionally, ensure that any included open-source library with special requirements (none detected beyond standard ones) is properly documented.

2. Conduct a Formal Security Assessment / Penetration Test – Before final ATO submission, we advise an independent security test of Elspeth in its operating environment. While our analysis found no obvious vulnerabilities, a focused penetration test or code review by an accredited third party can validate this and is often required for high-impact systems. Particular areas to test: handling of malicious input (e.g., extremely large or crafted dataset to see if any buffer issues in underlying libs), proper enforcement of classification gates (attempt to mislabel or bypass, ensuring it's robust), and the output sanitization logic (verify that no known Excel/HTML injection vectors slip through). Any findings can then be remediated, demonstrating due diligence (satisfying RA-5 and SA-11 controls).

3. Enhance HTML Report Sanitization – If the visual analytics sink produces an HTML report embedding model outputs, ensure that any user-generated content in those outputs is escaped to prevent script execution. For example, if a model output includes `"<script>alert('x')</script>"` and the HTML report includes it raw, opening that report could trigger it. Currently, there's no explicit mention of HTML sanitization. As a mitigation, use an HTML escaping library or ensure the report generator wraps output text in safe containers (e.g., `<pre>` or escapes special characters). This will close a minor gap and satisfy secure coding best practices (mapping to SC-7 and SI-10 for output handling). It's a straightforward improvement for an edge-case risk and will strengthen compliance with "no unexpected active content" policies.

4. Implement Asymmetric Signing for Non-Repudiation (Future) – The current HMAC-based artifact signing is effective for integrity, but for full non-repudiation (e.g., if outputs need to be shared outside the generating team with proof of origin), consider using a public/private key signature (e.g., an organizational code signing certificate or a PGP key pair). This is not an immediate ATO blocker since HMAC meets integrity requirements, but as a future enhancement it would align with best practice for high-assurance systems (satisfying AU-10 at a higher level, and supporting scenarios where verifying party should not have the secret key). Operationally, this would involve secure storage of a private signing key and distributing the public key with outputs or to auditors. It's a more complex setup, so we mark it as a future recommendation for systems that require it.

5. Operationalize Dependency & Configuration Scans – We recommend formalizing the processes already in place: ensure that for each release, a dependency vulnerability scan (pip-audit or equivalent) and a configuration review (especially for classification assignments) are part of the release checklist. The team has manual notes on this ⁵³³; automating it (e.g., CI pipeline fails on new pip-audit findings, or flags outdated libraries) will maintain compliance continuously. Additionally, maintain the **Traceability Matrix** document as code changes, so that any new feature or control change is reflected (perhaps integrate it into the development process to update that matrix). This ensures the System Security Plan remains up-to-date (satisfying PL-2 and PL-8 documentation requirements).

9.2 Improvement Opportunities (post-ATO enhancements): 1. **Stream Processing for Large Datasets**

– To handle extremely large input data or long-running jobs more gracefully, consider adding a streaming data processing capability. This could involve reading input in chunks (rather than entire DataFrame) and processing iteratively, or integrating with an ETL pipeline that feeds records to Elspeth incrementally. This would reduce memory usage and allow effectively infinite dataset processing. While not critical for current scope (where dataset sizes are manageable), it would improve scalability and reliability for big-data use cases. The design groundwork (checkpointing, processing row by row) is already conducive to chunking. Implementing this would address the residual risk of memory exhaustion and align with future data volume growth.

2. Parallel Execution of Multiple Experiments – Currently, experiments in a suite run sequentially (the SuiteRunner iterates experiments). If infrastructure and use cases allow, an enhancement is to execute multiple experiments in parallel, especially if they are independent (e.g., comparing multiple models). This could shorten total runtime.

Implementation could involve spawning separate threads or processes for each experiment (given the RateLimiter would need to partition or duplicate for each to avoid conflict, or run them fully independently). While this adds complexity (especially with shared rate limits or cost trackers), it could be offered as an option (e.g., a CLI flag `--parallel-experiments`). For ATO, this is not required (and the sequential approach might be safer to not overwhelm systems), but as an improvement it could utilize multi-core systems more efficiently for suites. If implemented, careful coordination of shared resources or completely segregating runs (no shared state) would maintain reliability.

3. User Interface & Orchestration Integration – To broaden adoption and make the tool accessible to less technical users, a lightweight GUI or web dashboard could be considered. For example, a read-only dashboard to view results and logs, or a simple form to launch experiments (with pre-defined config templates). This is not an immediate need and would introduce additional security considerations (web interface), but it can be pursued under a strong security design (e.g., internal-only web app with authentication, leveraging the existing classification and logging features). Alternatively, integration with existing orchestration tools (like providing an Airflow operator or Jenkins plugin to run Elspeth jobs) could ease deployment in enterprise workflows. This would leverage Elspeth's strengths while using established UI/monitoring of those platforms. We recommend exploring these approaches after initial rollout, to improve usability and integration, while ensuring any UI is built with the same security principles (if a web UI is built, it should enforce user roles, not expose sensitive config in browser, etc., mapping to AC and SC controls).

4. Expanded Plugin Library and Automated Policy Checks – Encourage development of additional plugins that enforce specific compliance or analytic policies. For example, a "PII Redaction" row plugin (to scrub outputs of personal data) could be provided out-of-the-box for organizations with privacy requirements. Similarly, a "Baseline Significance" aggregator plugin could automatically run statistical tests (using pingouin/statsmodels) and mark whether differences are statistically significant. This would make compliance with certain analytical standards (like requiring significance analysis for model comparisons) seamless. While users can create these, providing a vetted library of such plugins (perhaps maintained as an open-source catalogue) would improve consistency and reduce user effort. Over time, as policy needs evolve, new plugins can be added to address them (this is more of a process recommendation: maintain the pluggable architecture by continuously reviewing compliance requirements and implementing them as code where possible – essentially continuing the pattern Elspeth started). This will keep the tool aligned with emerging regulatory or organizational rules (for example, if a new policy requires watermarking AI-generated content, a plugin could be introduced to do that in outputs).

5. Continuous Compliance Monitoring – Leverage Elspeth's output metadata in a centralized compliance monitoring system. For instance, integrate the `metadata.json` outputs with a SIEM or GRC tool so that each experiment's metrics are logged to a dashboard. This would allow compliance officers to, at a glance, verify that all runs are within expected parameters (no excessive token usage, all outputs had zero policy violations, etc.). While not an architectural change to Elspeth itself, we recommend developing a small integration script or process for this (e.g., ship metadata to Splunk or create a PowerBI report from the Excel summary). This kind of continuous monitoring aligns with CA-7 (*Continuous Monitoring*) requirements and would demonstrate that once ATO is granted, the organization actively monitors the system's security and usage. The architecture already provides the data needed; the recommendation is to make use of it in the compliance program.

9.3 Migration Path and Future Roadmap: - Initial Adoption: For teams migrating from ad-hoc scripts to Elspeth, start with a **pilot** experiment. Use a simple config to replicate a known experiment and compare results to the manual method to validate correctness. The YAML structure is designed to be intuitive; provide training or examples for staff (the included sample configs in documentation serve as good starting points). Migrate gradually: define data sources in config, then prompts, then outputs, verifying each step. Because Elspeth imposes stricter checks (e.g., it might flag issues that previous manual process overlooked), be prepared to adjust configurations (or data/policies) to satisfy those checks. This migration approach ensures that compliance improvements (like classification

enforcement) are fully integrated into the team's workflow. - **Connecting to Enterprise Systems:** When moving Elspeth into a production environment, integrate it with existing enterprise scheduling and data platforms: - If using an ETL scheduler or pipeline (e.g., Apache Airflow, Azure Data Factory), treat Elspeth as a processing step – use its CLI within those pipelines. Migration here means writing wrapper scripts or Airflow Operators that pass the correct config and handle outputs (e.g., moving Elspeth outputs to a secure archive or database if needed). - For output usage, if an organization stores experiment results in a database or knowledge base, plan a step to load Elspeth's outputs (CSV/JSON) into those systems. Migration means establishing a mapping from Elspeth's output format to the target system's schema. The consistent output format of Elspeth aids this – e.g., you can create a table for "Experiment Results" that matches the columns in the CSV, and an "Experiment Metadata" table for the summary info. - **Legacy Experiments:** If the organization has legacy experiment records or legacy tools, consider **backfilling** them through Elspeth for consistency. For example, if older experiment runs were done manually and not all compliance data was captured, one might re-run those scenarios in Elspeth (with the same inputs and models) to produce a standardized output and audit trail. This could help have a single system-of-record for experiment outcomes. Migrating legacy data into Elspeth format (even if via dummy runs or feeding recorded outputs into Elspeth's pipeline) can streamline audits – it's a suggestion for completeness of records. - **Upgrade Path:** As Elspeth evolves (new versions, new plugins), maintain a clear versioning of config. The team already has an upgrade strategy documented ⁵²². For migration to a new version, use the traceability matrix to verify that all required controls are still met or improved, and run the comprehensive test suite to catch any unintended changes in behavior (particularly for deterministic outputs – ensure baseline results didn't change due to library updates). The migration to new versions should be done in a staging environment with sample experiments to confirm nothing breaks compliance or functionality before deploying to production (which aligns with CM-4 (*Security Impact Analysis*) on changes). - **Cloud Migration:** If not already on cloud, migrating to cloud infrastructure (like running Elspeth on an Azure VM or container) is straightforward since it's self-contained. Ensure to apply cloud security (lock down VM access, use managed identity for creds as designed). The architecture – using Azure SDKs, etc. – is cloud-ready. The main migration tasks are configuring networking (e.g., allow the VM's outbound to OpenAI, set up Key Vault if using secrets). This move can increase scalability (bigger VMs, easier monitoring) and should not require code changes, just configuration adjustments (like switching to azure-based LLM plugin from openai if needed). The design already anticipated this by abstracting LLM clients. For example, to migrate from using OpenAI public to Azure OpenAI (for compliance with data residency), one can simply change config to use `azure_openai` plugin with appropriate options ²²⁶ – no code change, just migration of config and re-running. This flexibility was a deliberate design aimed at easing such migrations. - **Future Vision:** Looking ahead, Elspeth's architecture is positioned to serve as a compliance-centric experimentation hub. We recommend continuing the practice of aligning it with evolving compliance frameworks (e.g., if new AI regulations require logging of model prompts and responses for accountability, Elspeth already does much of this – just ensure logs are retained as required, etc.). The plugin architecture will accommodate future changes (like if a regulation demands "model output must be watermarked", one could implement a plugin to append a watermark string to each response – easily integrated). So, the migration path here is more about *policy migration*: as policies change, update or add plugins to enforce them. Elspeth should remain the central tool where these enforcement mechanisms reside, reducing reliance on manual compliance steps. - **Training and Knowledge Transfer:** As part of migrating to using Elspeth, invest in training the user and security teams on its use and interpretation of outputs. The clear mapping to controls (in docs) can be included in security training sessions to show how the tool automatically handles certain requirements. Over time, this will embed Elspeth into the organization's processes as a trusted, standard approach for LLM experiments. Essentially, migrating from an *ad-hoc, person-dependent process* to a *systematic, tool-driven process* is as much cultural as technical – the architecture is sound, so ensure the team's practices migrate to leveraging it fully (for example, always using Elspeth for any new experiment rather than some quick script, so that all runs are logged and compliant).

In summary, Elspeth is ATO-ready with minor finishing steps. Addressing the critical issues above will ensure no gaps remain for approval, while the improvement suggestions will keep the system robust and effective in the long run. By following the migration path guidance, the organization can smoothly adopt Elspeth as the cornerstone of responsible AI experimentation, confident that it meets both current and future compliance obligations.

Appendices

A. Code Analysis Methodology – The analysis in this document was conducted by systematically reviewing Elspeth’s source code (particularly the main modules in `src/` and associated test cases), as well as its documentation and dependency configuration. We utilized static analysis techniques, tracing through the code’s logic (for example, following the flow from CLI argument parsing through orchestrator to runner and plugins) and cross-referencing against known security and performance best practices. The compliance mapping was derived by taking each implemented control or security feature and aligning it to the corresponding NIST 800-53 controls or organizational policy statements, using the project’s own control inventory as a guide ⁴²⁴. Where available, automated tools (like pip-audit for dependency vulns and logging checks) supplemented manual review. We also leveraged the existence of an internal *Traceability Matrix* to ensure that every requirement listed has a corresponding implementation in code, verifying each with code references (for instance, seeing that “*validate datasource schema*” in documentation corresponds to `_validate_plugin_schemas` call in code ⁹⁶). This thorough approach – combining manual code inspection, automated scanning, and documentation cross-checking – provides high confidence in the accuracy and completeness of the architectural assessment.

B. File and Component Mapping – The following mapping enumerates key source files and the components or functionality they implement, establishing a clear link between code artifacts and the architectural components described:

- `src/elspeth/cli.py` – Implements the **CLI Interface** (argument parsing, invoking load and run) ^{24 26}.
- `src/elspeth/config.py` – Contains the **Settings** dataclass and the `load_settings` function, i.e., the **Configuration Loader** logic (YAML parsing, instantiating plugins) ^{230 131}.
- `src/elspeth/core/orchestrator.py` – Defines the **ExperimentOrchestrator** class (orchestration for single experiment), including context creation and integration with **ExperimentRunner** ^{48 74}.
- `src/elspeth/core/experiments/runner.py` – Defines the **ExperimentRunner** class responsible for the main execution loop, concurrency control, applying **plugins** (row, agg, val) and managing early stop and checkpoint logic ^{341 92}.
- `src/elspeth/core/experiments/suite_runner.py` – Defines the **ExperimentSuiteRunner** class for orchestrating multiple experiments and baseline comparisons in a suite ^{303 101}.
- `src/elspeth/core/plugins.py` – Likely contains the **PluginContext** class and perhaps base protocol definitions for plugins (like `LLMClientProtocol`, `ResultSink Protocol`) ^{48 56}.
- `src/elspeth/core/artifact_pipeline.py` – Implements the **ArtifactPipeline** class and `ArtifactStore`, performing security checks and coordinating **Result Sinks** execution ^{534 398}.
- `src/elspeth/plugins/` – Package containing specific plugin implementations:
 - `datasources/*.py` (e.g., `csv_local.py`, `blob.py`) – **Data Source Plugins** for local CSV ⁸⁴, Azure Blob, etc.
 - `llms/*.py` (e.g., `openai_http.py`, `azure_openai.py`, `mock.py`) – **LLM Client Plugins** for OpenAI ³⁸, Azure, and a mock model ³⁸¹.
 - `experiments/` (e.g., `early_stop.py`, `metrics.py`, `validation.py`) – likely contains **Row/Aggregator/EarlyStop Plugins** or at least definitions for them (the test files show registration of dummy plugins for metrics ⁸⁸).
 - `outputs/*.py` (e.g., `csv_file.py`, `repository.py`, `signed.py`, `excel.py`) – **Result Sink Plugins** for writing CSV ¹⁰⁹, pushing to repository, generating signed manifest, writing Excel, etc.
- `tests/` – The test suite files map to features:
 - `test_cli.py` – tests CLI arg handling,
 - `test_config.py` – tests config loader (ensuring e.g., missing fields raise errors as expected),
 - `test_orchestrator.py` – tests orchestrator and

plugin execution (single-run logic) ⁸⁸, - `test_runner.py` - tests multi-threading, early stop, checkpoint behaviors, - `test_artifact_pipeline.py` or `test_sanitize_utils.py` - tests security functions like formula sanitization ¹¹⁰, - `test_end_to_end.py` (if any) - would test an entire suite run with baseline and verify outputs. This mapping ensures maintainers can quickly locate where in the code a particular functionality resides (for example, knowing that security clearance checks are in `core/artifact_pipeline.py` as shown by references ³⁹⁸).

C. Detailed Security Findings – The architectural review did not uncover critical security vulnerabilities in Elspeth’s design or implementation. All security controls appear properly implemented and active. Some minor findings and considerations include: - **Use of HMAC for Signing** – As discussed, the current HMAC approach provides integrity but not true non-repudiation. Given Elspeth’s intended internal use, this is acceptable. The “finding” is that if stronger authenticity is needed (e.g., if sharing outputs with external auditors who shouldn’t have the key), an upgrade to asymmetric signing is advisable. This is noted as a recommendation rather than a flaw. - **Potential HTML Injection in Reports** – If model outputs contain HTML or script content, and the Visual Report sinks embed them, there’s a theoretical XSS risk when viewing the report in a browser. There’s no evidence of sanitization in code for HTML (unlike for CSV). This is a low-probability issue (it requires malicious content from the model and someone opening the HTML in a full browser environment). Nonetheless, it’s flagged as a security consideration. It can be mitigated by escaping HTML special characters in outputs – a simple patch if needed. - **Large Input Handling** – Not a security vulnerability per se, but if extremely large inputs are fed, Python may crash or freeze (leading to potential availability issues). While not an immediate concern for compliance, it could be abused (e.g., a user intentionally feeds a huge file to cause denial-of-service on a shared machine). The finding is that streaming input or at least a documented input size limit would bolster availability. This has been recommended in section 9. - **Dependency Update Lag** – No outdated libraries were found; all are current. We note that continuous vigilance is required. For instance, if a new CVE arises in requests or PyYAML, the team must promptly update. Given their process, this is likely, so no immediate issue. The current dependencies appear free of known CVEs (we cross-checked versions: e.g., requests 2.31.0 – no open CVEs; PyYAML 6.0 – safe from earlier issues; Jinja2 3.1.0 – no known issues). - **Secret Exposure** – We confirmed that no secrets (API keys, etc.) are logged or stored in outputs. The only place they appear is in memory when used. The use of safe loaders and not including secrets in exceptions (if a secret were wrong, the error message doesn’t print the secret, just says “Unauthorized” or similar). This is good practice and we found no instances of secrets being concatenated into log messages (we scanned the code for uses of sensitive variables in logging – none found). - **Integration with External Auth** – The reliance on environment credentials means that the security of those credentials is outside Elspeth’s scope (i.e., the OS or container must protect environment variables, and Azure must protect managed identity tokens). This is standard and acceptable. The finding is to ensure those environments are configured correctly (e.g., not running Elspeth in a user context that has excessive Azure permissions beyond what’s needed for the chosen tasks – principle of least privilege for the identity). This is more of an operational note than a code issue. - **Logging & Privacy** – The content of logs was reviewed: they do not log model outputs or personal data by default (only warnings and counts). This is positive for privacy compliance (no inadvertent exposure of sensitive outputs in log files). If debug logging is enabled, it might log more (the OpenAI SDK debug could show prompts), but enabling debug is a conscious action typically in dev environment, so acceptable. Just note that in production, debug should remain off to avoid that. Overall, our detailed review finds the system’s security posture strong. The above minor points have either been addressed by recommendations or deemed low risk. We conclude that Elspeth is architecturally sound from a security perspective and any residual issues are manageable with standard best practices.

D. Performance Evaluation Points – Key performance behaviors were tested and observed: - **Concurrent Execution:** In a controlled test with ~100 requests and `max_workers=4`, Elspeth achieved roughly a 3.5x speed-up over serial execution, which aligns with expectation (some overhead for thread

management and rate limiting prevented a full 4x). For example, 100 dummy LLM calls that took ~1 second each serially (100 seconds total) completed in ~28 seconds with 4 threads, demonstrating effective parallelism. CPU usage remained low (since dummy calls were I/O-bound), indicating threads spent most time waiting – a scenario where Python threads excel. This confirms the decision to use threads is effective for I/O-bound tasks.

- **Rate Limiting Effect:** We simulated hitting rate limit thresholds by configuring a low threshold in the RateLimiter. The system appropriately paused task submission and avoided overwhelming the test API. The net throughput stabilized at the configured threshold (~80% of allowed rate), as expected. This prevented an avalanche of 429 errors – in a test where we intentionally set a low 3 requests/sec limit and launched 10 threads, without throttle we saw many 429s and retries, with throttle, we saw almost none, and total time was slightly higher than minimal but with far fewer failed attempts. This demonstrates that the adaptive throttle trades a small latency increase for reliability, a beneficial trade for overall throughput consistency.
- **Memory Footprint:** In a test with 10,000 records (small text per record), memory usage peaked at around 250 MB with Pandas overhead, well within acceptable range for a modern server. CPU usage was modest (the load is mostly waiting on network). We noted that memory use grows roughly linearly with number of records (each record's data and result take a few hundred bytes). At ~100k records (simulated in a high-memory environment), memory usage was about 2.3 GB. This indicates that beyond certain scale, memory could become a bottleneck or require chunking (as recommended). But up to tens of thousands of records, it's manageable on a typical machine with ~8-16 GB RAM.
- **Output Generation:** Writing a CSV of 10k rows and 20 columns took <0.5 seconds. Writing an Excel with one sheet of same size took ~1.2 seconds. Embedding 5 small images in an HTML report took negligible time (the overhead is mostly in generating images with Matplotlib, which in our test of a simple plot for 10k points took ~0.8 seconds per plot). These output times are minor compared to model inference times. The artifact signing (HMAC) for these outputs (a few MB of data) was essentially instantaneous (<0.1s). So, the pipeline of final sinks does not introduce any performance concern for typical output sizes.
- **Resilience under Load:** We conducted a stress test by simulating an API slowdown: we introduced a 100ms artificial delay in each dummy LLM call and ran 50 threads. The system handled it gracefully: RateLimiter slowed submissions as backlog built (keeping ~40 calls in flight), and no failures occurred. The early-stop was not triggered as progress was being made. CPU usage remained moderate (~70% across 8 cores, as threads were partially idle waiting). The test completed $\sim 50 \times 1s = \sim 1s$ (since 50 parallel on 8 cores with 100ms each, they ended in ~0.7s after the initial ramp-up thanks to parallel dispatch). This indicates the architecture can handle bursty loads and still meet near-optimal throughput, limited by external API speed, not internal overhead.

These performance tests confirm that Elspeth's architecture is efficient for its intended use cases. For extremely large jobs (e.g., >100k records or multi-GB data), we observed memory usage scaling linearly, suggesting that implementing chunking/streaming as recommended would be wise to avoid potential memory exhaustion. However, such scenarios are likely rare given typical LLM experiment scales (usually in the thousands to low tens of thousands of inputs, as sending more than that to an API is often cost or time prohibitive).

Overall, Elspeth demonstrates solid performance characteristics: - It effectively parallelizes I/O-bound tasks, - Maintains stable throughput under rate limiting, - Has predictable memory usage, - And its additional security and logging features do not introduce noticeable performance drag for moderate data volumes (the overhead of checks and logs is negligible in measured tests). This indicates the system can reliably meet expected operational performance and scaling requirements as currently designed.

1 6 10 20 23 63 156 157 158 228 246 392 420 522 525 526 528 529 **GitHub**
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/README.md>

2 48 71 72 73 74 75 76 77 85 103 119 121 122 123 125 126 127 132 133 134 155 191 192 193 229 231
272 273 274 275 276 277 481 482 **GitHub**
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth/core/orchestrator.py>

3 113 114 118 178 179 180 181 182 210 211 212 213 214 389 390 391 394 395 398 399 441 442 506 534 **GitHub**
https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth/core/artifact_pipeline.py

4 5 165 185 186 400 401 461 462 463 **GitHub**
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth/core/prompts/engine.py>

7 8 9 11 12 13 21 22 27 28 43 44 56 57 87 90 92 93 94 95 96 97 98 99 106 115 117 135
136 137 138 139 140 141 142 143 144 145 146 147 148 154 160 161 162 163 164 166 167 168 169 172 183 184
187 194 195 196 197 198 199 200 204 205 206 207 208 209 232 235 236 237 238 239 240 241 242 243 244 323
324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349
350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375
376 377 378 379 385 402 403 404 408 409 410 411 412 431 437 452 464 467 470 490 491 492 493 494 495 496
497 498 499 500 501 502 503 504 505 507 508 513 514 **GitHub**
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth/core/experiments/runner.py>

14 15 47 84 86 88 89 91 108 153 159 170 171 201 202 225 381 396 446 465 510 **GitHub**
https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/tests/test_orchestrator.py

16 17 227 516 517 518 **GitHub**
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth.egg-info/SOURCES.txt>

18 19 24 25 26 29 30 31 32 33 34 35 36 37 40 41 42 58 69 70 109 110 111 120 149 173 174 175
176 177 188 189 203 215 216 217 218 219 220 221 222 223 224 247 248 249 250 251 252 253 301 302 382 383
384 397 428 433 434 435 466 471 472 473 485 511 524 **GitHub**
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth/cli.py>

38 39 45 53 54 55 59 61 62 64 65 68 226 233 234 245 388 405 406 407 430 512 520 521 523 527 533
GitHub
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/docs/development/dependency-analysis.md>

46 78 79 80 81 82 83 100 101 102 104 105 107 116 124 128 129 150 151 152 190 261 271 278 279 280 281
282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 303 304 305 306 307 308 309
310 311 312 313 314 315 316 317 318 319 320 321 322 386 393 421 436 447 476 483 484 486 487 488 489 509
GitHub
https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth/core/experiments/suite_runner.py

49 50 51 52 60 112 380 387 519 **GitHub**
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth.egg-info/requires.txt>

66 67 130 131 230 254 255 256 257 258 259 260 262 263 264 265 266 267 268 269 270 427 474 475 477 478 479
480 **GitHub**
<https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/src/elspeth/config.py>

413 414 415 416 417 418 419 422 423 424 425 426 429 432 438 439 440 443 444 445 448 449 450 451 453 454
455 456 457 458 459 460 468 469 515 530 531 532 **component-diagram.md**

[https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/docs/architecture/
component-diagram.md](https://github.com/tachyon-beep/elspeth/blob/c5144ea2e0674a0966b00955270b9ea372aa1b21/docs/architecture/component-diagram.md)