

Solution Overview

This project implements a **modular, efficient Log Analyzer** in Python, designed to process large sets of log files, apply flexible event-based filters, and output structured, human-readable results. The tool is intended for use in environments where logs are generated in high volume and need to be filtered and summarized according to dynamic criteria, such as in system monitoring, telemetry, or device management scenarios.

Code Structure and Flow

1. Entry Point and CLI

- The script is designed to be run from the command line, using Python's argparse module to provide a clear and user-friendly interface.
- Required arguments:
 - --log-dir: Directory containing log files (.log or .log.gz).
 - --events-file: Path to a configuration file specifying event filters.
 - Optional arguments: --from and --to for timestamp-based filtering.

2. Event Filter Parsing

- The parse_events_file function reads the events configuration file, which defines filters in a simple, extensible syntax (event type, optional log level, optional regex pattern, and a count flag).
- Each filter is represented by an EventFilter object, encapsulating its logic and description.

3. Log File Discovery and Reading

- The log_files_in_dir generator yields all valid log files in the specified directory, ensuring only files (not directories) are processed.
- The log_lines_from_file generator transparently supports both plain text and gzip-compressed logs, allowing the tool to handle large or archived datasets efficiently.

4. Log Line Parsing

- Each log line is parsed using a regular expression to extract the timestamp, log level, event type, and message.
- Malformed lines are skipped, and errors are logged for debugging.

5. Filtering and Analysis

- The `filter_logs` generator streams log entries, applying optional timestamp filters to avoid loading unnecessary data into memory.
- The `analyze_logs` function iterates through all logs and applies each event filter, collecting matches in a memory-efficient way (streaming, not batch loading).

6. Result Formatting

- The `format_results` function produces a structured, readable output for each filter:
 - If the filter is a count, it reports the number of matches.
 - Otherwise, it lists all matching log lines.
 - Separator lines are used for clarity between filters.

7. Error Handling and Robustness

- The tool uses Python's logging module for error and warning messages, ensuring that issues in input files or configuration are reported clearly.
 - All file and argument errors are handled gracefully, with user-friendly messages.
-

Design Considerations

1. Modularity and Extensibility

- The code is organized into small, single-responsibility functions and classes, making it easy to extend (e.g., adding new filter types or output formats).
- The `EventFilter` class encapsulates all logic related to a single filter, supporting future enhancements (such as more complex matching or aggregation).

2. Performance and Scalability

- The use of generators (`yield`) throughout the code ensures that the tool can process very large log datasets without excessive memory usage.
- Support for `.gz` files allows for efficient storage and processing of archived logs.

3. Usability and UX

- The CLI interface is clear and provides helpful error messages.
- Output is structured for both human readability and potential downstream processing.

4. Robustness and Testing

Yair Dayan – Log Analyzer

- The code is defensive against malformed input, missing files, and invalid configuration.
- Comprehensive unit tests are provided to verify correct behavior across a range of scenarios, including edge cases and error conditions.

5. Maintainability

- The code is well-documented, with docstrings for all major functions and classes.
 - Logging and error handling are consistent, aiding future debugging and maintenance.
-

Summary

This Log Analyzer is a robust, efficient, and extensible tool for extracting structured insights from large log datasets. Its modular design, memory efficiency, and clear CLI make it suitable for both ad-hoc analysis and integration into larger monitoring or automation pipelines. The codebase is well-tested and documented, ensuring reliability and ease of future development.