



# White Agent Report: LangGraph-Based Documentation Generation

## Abstract

We present a LangGraph-based documentation generation agent for AEO-Bench, a benchmark evaluating AI agents' ability to transform code into comprehensive documentation. Our agent employs an explicit state machine with three specialized nodes—PLANNER, EXPLORER, and GENERATOR—each with focused prompts optimized for its specific task. Evaluated on 3 test cases, the white agent achieves 235/300 (78.3%) total score compared to the baseline's 93/300 (31.0%), while completing 40% faster. This separation of concerns produces better documentation than monolithic conversation-based approaches by avoiding context overload.

## Benchmark: AEO-Bench

AEO-Bench evaluates project-level documentation generation, filling a gap left by benchmarks like SWE-Bench (Jimenez et al., 2024) which test coding ability, and DocPrompting (Zhou et al., 2022) which focuses on function-level docstrings. The benchmark uses 6 repositories—3 synthetic and 3 real GitHub repos—where agents explore code using `list_directory` and `read_file` tools before submitting documentation via `respond`.

Evaluation follows a 4-tier rubric totaling 100 points: Tier 1 Structural (15 pts) checks valid JSON and README presence; Tier 2 Sections (25 pts) detects Installation, Usage, and Examples via keyword matching; Tier 3 Accuracy (30 pts) uses an LLM judge to verify purpose, dependencies, and commands; Tier 4 Quality (30 pts) assesses clarity, completeness, and formatting.

## White Agent Architecture

The white agent uses a LangGraph state machine with three nodes that execute in sequence. The PLANNER node recursively explores the directory structure, identifying all files and subdirectories before creating a prioritized reading plan. The EXPLORER node then reads files

deterministically from this plan without requiring additional LLM calls—a key efficiency gain. Finally, the GENERATOR node synthesizes all accumulated file contents into a comprehensive README and [schema.org](#) metadata.

State tracking enables this modular approach: the agent maintains fields for phase (planning/exploring/generating), discovered files and directories, exploration history, file contents, and the current action to return. This explicit state allows systematic file reading and ensures comprehensive codebase coverage across request-response cycles.

The design's key advantages include explicit planning before exploration, recursive directory handling for nested packages, deterministic file reading that eliminates LLM overhead during exploration, and modular nodes that can be modified independently.

## Experiments

We compared the white agent against a baseline single-prompt agent on 3 test cases. The white agent scored 235/300 total (78.3% average) while the baseline scored 93/300 total (31.0% average), with the white agent completing in 49.2s versus the baseline's 81.7s. On individual cases, the white agent achieved 87/100 on art\_github, 85/100 on countdown\_timer, and 63/100 on dotenv\_github. The baseline scored 93/100 on art\_github but failed completely on the other two cases due to JSON parsing errors.

The white agent consistently achieves strong Tier 1 (15/15) and Tier 2 (25/25) scores with good Tier 3 (22-30/30) and Tier 4 (25/30) performance. The baseline struggles with output formatting, frequently producing malformed JSON that fails evaluation—a consequence of accumulated conversation history becoming unwieldy and causing the model to lose track of formatting requirements.

Token efficiency favors the white agent significantly. While the baseline accumulates growing conversation history with each step, eventually causing context overload, the white agent uses fresh specialized prompts for each node. EXPLORER makes no LLM calls at all, PLANNER makes one call per directory explored, and GENERATOR makes a single call for final output. This results in 40% faster completion with higher quality output.