
RPG: A REPOSITORY PLANNING GRAPH FOR UNIFIED AND SCALABLE CODEBASE GENERATION

Jane Luo^{1*,‡} Xin Zhang^{1*†} Steven Liu^{1,‡} Jie Wu^{1,2,‡} Jianfeng Liu¹ Yiming Huang³ Yangyu Huang¹
 Chengyu Yin¹ Ying Xin¹ Yuefeng Zhan¹ Hao Sun¹ Qi Chen¹ Scarlett Li¹ Mao Yang¹

¹ Microsoft ² Tsinghua University ³ University of California, San Diego

October 13, 2025

ABSTRACT

Large language models excel at generating individual functions or single files of code, yet generating complete repositories from scratch remains a fundamental challenge. This capability is key to building coherent software systems from high-level specifications and realizing the full potential of automated code generation. The process requires planning at two levels: deciding what features and modules to build (proposal stage) and defining their implementation details (implementation stage). Current approaches rely on natural language planning, which often produces unclear specifications, misaligned components, and brittle designs due to its inherent ambiguity and lack of structure. To address these limitations, we introduce the Repository Planning Graph (RPG), a structured representation that encodes capabilities, file structures, data flows, and functions in a unified graph. By replacing free-form natural language with an explicit blueprint, RPG enables consistent long-horizon planning for repository generation. Building on RPG, we develop ZeroRepo, a graph-driven framework that operates in three stages: proposal-level planning, implementation-level construction, and graph-guided code generation with test validation. To evaluate, we construct RepoCraft, a benchmark of six real-world projects with 1,052 tasks. On RepoCraft, ZeroRepo produces nearly 36K Code Lines and 445K Code Tokens, on average 3.9× larger than the strongest baseline (Claude Code), and 68× larger than other baselines. It achieves 81.5% coverage and 69.7% test accuracy, improving over Claude Code by 27.3 and 35.8 points. Further analysis shows that RPG models complex dependencies, enables more sophisticated planning through near-linear scaling, and improves agent understanding of repositories, thus accelerating localization.

1 Introduction

Recent large language models (LLMs) have shown strong performance on function-level and file-level code generation, reliably producing functions and files from natural language descriptions [1–4]. However, scaling this capability from functions and files to generate large-scale software repositories from scratch remains a fundamental challenge. The core difficulty is bridging the gap between high-level user intent and the repository’s intricate network of files, classes, and dependencies[5, 6]. Successfully navigating this gap necessitates a process of progressive planning, which naturally decomposes into two complementary phases: **proposal-level planning**, which determines *what to build* by defining the functional scope and key capabilities, and **implementation-level planning**, which determines *how to build* it by specifying the file structure, interfaces, dependencies, and data flows.

Prior work has explored this challenge through three paradigms. Distributed planning frameworks (e.g., MetaGPT [7], ChatDev [8]) assign specialized roles such as manager, architect, and engineer to negotiate between requirements and implementations. Workflow-based systems (e.g., Paper2Code [9], AutoP2C [10]) follow fixed pipelines that first build

* Equal contribution, Contact: janeluo1210@163.com

† Corresponding author: xinzhang3@microsoft.com

‡ This work is done during their internships at Microsoft.

architectural skeletons before filling in details. Iterative terminal agents (e.g., OpenHands [11], Claude Code [12], Gemini CLI [13]) externalize intermediate plans, often in markdown, and refine them step by step. Despite their differences, these approaches share a dependency: natural language as the intermediate medium for planning.

While natural language remains a flexible and human-readable medium, it can often be less efficient for large-scale repository generation. Its inherent ambiguity may blur distinctions between intent and constraints [14], its lack of explicit hierarchy makes dependency tracking particularly difficult [15], and static plans may gradually degrade over long horizons without adaptive adjustment [16]. When extended to automatic repository generation, these limitations can more easily lead to unstable proposal-level planning, where functionalities are sometimes incomplete, overlapping, or unevenly scoped [17], and fragmented implementation-level planning, where plans drift across iterations, introducing inconsistencies in dependencies, data flows, and modular boundaries [18, 19].

To address these limitations, we introduce the Repository Planning Graph (RPG), a persistent and evolvable representation that unifies proposal and implementation planning for repository generation. RPG encodes functional goals and designs in a single graph: nodes capture hierarchical capabilities with files, classes, and functions, while edges specify semantic relations and data flows. By replacing free-form language with a structured medium, RPG provides a compact, interpretable basis for consistent long-horizon planning. Building on this representation, we develop ZeroRepo, a graph-driven framework for controllable repository generation. Given a user specification, ZeroRepo proceeds in three stages: (1) **Proposal-Level Construction**, which organizes and refines requirements into a functional graph via a large-scale feature tree; (2) **Implementation-Level Construction**, which expands this graph into the full RPG by encoding file skeletons, interfaces, and flows; and (3) **Graph-Guided Code Generation**, which traverses the RPG in topological order with test-driven development, guided localization, and iterative editing.

To evaluate agents’ planning ability in repository generation, we construct **RepoCraft**, a benchmark of six projects with 1,052 tasks. On RepoCraft, ZeroRepo attains 81.5% functional coverage and a 69.7% pass rate, exceeding the strongest baseline (Claude Code) by 27.3 and 35.8 points, while producing repositories with 36K Lines of Code and 445K Code Tokens, about 3.9× larger than Claude Code and 68× larger than other baselines. Further analysis shows that **Repository Planning Graph (RPG)** captures complex dependencies, including inter-module data flows and function-level relations. It enables near-linear scaling of functionality and code size, supporting complex planning and providing a foundation for large-scale repositories and long-horizon development. As a global representation, RPG enhances agents’ repository understanding and accelerates localization.

Our Contributions are list below:

1. We introduce the Repository Planning Graph (RPG), a unified representation integrating proposal- and implementation-level planning, encoding functionality, file structures, data flows, and function designs.
2. We develop ZeroRepo, a graph-driven framework that constructs RPG through proposal- and implementation-level planning, and generates code with test validation.
3. To evaluate agent planning ability in repository generation, we build **RepoCraft**, a benchmark of 6 projects with 1,052 tasks assessing coverage, accuracy, and code scale.
4. On RepoCraft, ZeroRepo achieves strong improvements over baselines, reaching 81.5% functional coverage and nearly 69.7% test accuracy, while producing repositories 3.9× larger than the strongest baseline. Further analysis shows that RPG captures complex dependencies, enables more sophisticated planning through near-linear scaling, and enhances agents’ repository understanding, thereby accelerating localization.

2 Related Work

LLM-based Code Generation SOTA models (e.g., GPT-4o [20], Claude 4 [21], Gemini [22], DeepSeek-R1 [23]) excel at diverse SWE tasks, including code completion, test generation [24, 25], refactoring [26], and program repair [27]. Instruction-tuned variants (e.g., Qwen-Coder [28], EpiCoder [2]) further improve reliability. These advances establish strong function-level performance, laying the foundation for progress toward broader software engineering tasks.

Agents for Repository-Level Generation Agent frameworks embed LLMs in planning–coding loops [29, 30]. Multi-agent systems (e.g., ChatDev [8], MetaGPT [7]) assign roles under fixed procedures, while workflow systems (e.g., Paper2Code [9], AutoP2C [10]) decompose synthesis into phases. Industrial systems (e.g., Codex [31], Gemini Cli [32] Claude Code [12]) extend these ideas to multi-file projects. However, most rely on ad-hoc natural language plans without persistent structural representations, often leading to fragmented implementations. By contrast, ZeroRepo employs a graph-guided abstraction that enforces structured planning and implementation guide.

3 Repository Planning Graph Construction

To address the ambiguity of natural language plans, we propose the **Repository Planning Graph (RPG)**, a structured representation that encodes repository functionality and implementation logic as nodes and edges. Building on RPG, we

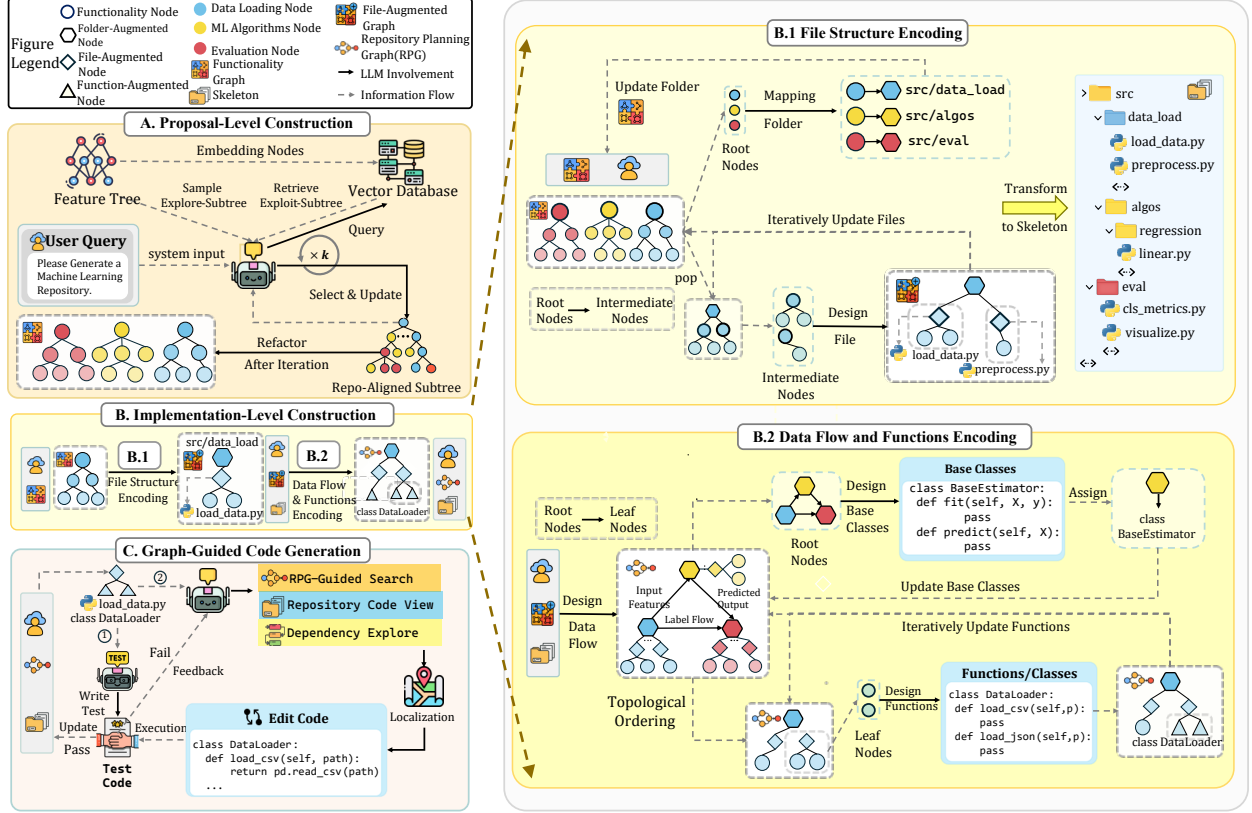


Figure 1: The ZeroRepo pipeline for repository generation. (A) Proposal-level construction translates specifications into a functionality graph. (B) Implementation-level construction refines it via (B1) file-structure encoding into a file-augmented graph and (B2) data-flow/function encoding into the final Repository Planning Graph (RPG). (C) Graph-guided code generation traverses RPG in topological order to produce a stable repository.

develop ZeroRepo, a framework for repository generation from scratch. This section first introduces the structure of RPG (§3.1), and then explains how ZeroRepo constructs it through proposal-level planning (§3.2) and implementation-level refinement (§3.3). The overall pipeline is shown in Figure 1(A-B).

3.1 Repository Planning Graph Structure

As shown in Figure 2, RPG provides a unified representation for repository planning by encoding functionality and implementation in a structured graph, rather than unstable natural language descriptions. Its nodes carry dual semantics. At the functional level, they represent progressively refined capabilities: high-level modules (e.g., algorithms, evaluation) decompose into mid-level components and ultimately into leaf nodes corresponding to concrete algorithms. At the structural level, this hierarchy naturally mirrors repository organization: root nodes typically align with file regions, intermediate nodes with files, and leaf nodes with functions or classes, thereby unifying functional decomposition with code structure.

Beyond the hierarchical nodes, edges in RPG capture execution dependencies across granularity. Inter-module edges (black arrows in Figure 2) encode data flows between modules, such as outputs from Data Loading feeding into ML Algorithms and then into Evaluation. Intra-module edges (gray dashed arrows) capture file-level orderings; for instance, `load_data.py` precedes `preprocess.py`, with

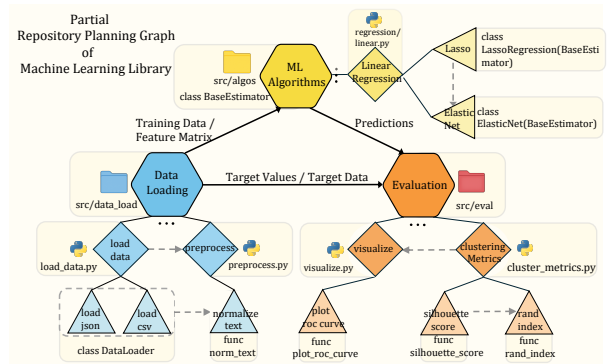


Figure 2: Example of a Repository Planning Graph: Solid lines show hierarchy, black arrows for inter-module flows, dashed gray arrows for intra-module order.

outputs propagated to preprocessing. Collectively, these edges impose a topological order that aligns functional decomposition with code organization, ensuring coherence between global execution semantics and local implementation.

3.2 Proposal-Level Construction

At the proposal level, the aim is to translate high-level user specifications into a coherent functionality graph. This involves three steps: grounding functionalities in a large-scale Feature Tree, selecting a repository-aligned subtree via explore–exploit search, and refactoring it into the graph. The full algorithm is detailed in Appendix A.1.

A Global Tree as Knowledge Base LLMs alone provide unstable and biased capability enumeration, often yielding incomplete coverage [33, 34]. To stabilize planning, we ground functionality selection in the EpiCoder Feature Tree [2], a large-scale ontology of over 1.5M software capabilities. Its broad coverage and hierarchy serve as a structured prior, mitigating randomness and bias while ensuring repository functionalities are systematically captured and diverse. For efficient retrieval, each feature node is embedded into a vector representation, with its full hierarchical path stored as metadata in a vector database. This design preserves both semantic similarity and structural context, enabling precise and scalable functionality grounding. Detailed statistics of the Feature Tree are provided in Appendix A.2.

Explore–Exploit Subtree Selection Using the Feature Tree as a structured knowledge base, the first step is to construct a **repo-aligned subtree** tailored to the user’s goal. Exhaustive enumeration is infeasible at the 1.5M scale, so ZeroRepo incrementally expands the subtree via an explore–exploit strategy. (1) **Exploitation** ensures precision: we retrieve top- k feature paths most aligned with the user goal and augment them with keywords suggested by LLM queries. (2) **Exploration** ensures diversity: we deliberately expand into unvisited regions of the ontology to capture less obvious but relevant functionalities. Candidates from both strategies are filtered by the LLM and integrated into the evolving subtree, yielding a balanced and comprehensive foundation for downstream planning.

Refactoring by Goal Alignment The repo-aligned subtree, though capturing relevant functionalities, still inherits the generic organization of the global ontology. To align it with the user’s repository goal, we refactor feature placements into a modular **functionality graph**. The LLM partitions functionalities into cohesive modules following software engineering principles of cohesion and coupling. For instance, in a machine learning library, metrics such as `silhouette_score` are reorganized under an evaluation module rather than within clustering algorithms. The resulting graph establishes clear functional boundaries, encoding proposal-level planning directly into the representation.

3.3 Implementation-Level Construction

After proposal-level construction establishes the multi-level functional plan, the graph is further enriched with implementation details, culminating in the complete Repository Planning Graph (RPG) at this stage. The process includes encoding the repository’s file structure, modeling inter-module data flows and intra-module orderings, and specifying concrete functions and interfaces.

3.3.1 File Structure Encoding

While proposal-level planning defines modular functionalities, it remains abstract and detached from implementation. To bridge this gap, the graph is extended with folder and file layouts, instantiating a repository skeleton that maps functional modules into executable structures, resulting in a file-augmented graph.

Folder-Level Encoding Proposal-level planning partitions functionalities into modular subgraphs, but this abstraction does not yet define the repository’s structural layout. To bridge the gap, we enrich root nodes with folder-level specifications, assigning each subgraph a dedicated directory namespace (e.g., `algorithms/`, `evaluation/`). This encoding couples semantic modularity with explicit structural separation, ensuring that descendant functionalities inherit a consistent namespace and that the repository skeleton align with high-level capability decomposition.

File-Level Encoding Once folder regions are encoded at the root nodes, the graph is enriched by assigning files to intermediate nodes. This step specifies how functionalities within a module are grouped into executable files. For example, preprocessing utilities are consolidated into `preprocess.py`, while models such as linear regression and its variants are grouped into `linear_models.py`. By embedding file structure in the graph, we preserve semantic cohesion, reduce cross-file coupling, and obtain a file-augmented graph that anchors downstream design.

3.3.2 Data Flow and Functions Encoding

After obtaining the file-augmented graph, this stage finalizes the full Repository Planning Graph (RPG) by assigning executable roles to leaf nodes. To ensure coherence across modules and functions, we first incorporate inter- and intra-module data flows as input–output constraints, then abstract shared structures as design anchors, and finally refine leaf nodes into concrete functions or classes.

Data-Flow Encoding To ground interface design in execution semantics, the graph is augmented with data-flow edges that capture inter- and intra-module relations. At the global level, as shown in Figure 2, typed input–output flows connect subgraph roots; for example, a data-loading module may provide an array of training data to downstream algorithms. At the local level, files within a module are planned in a specific order, ensuring coherent and dependency-aware implementation. These flows impose a hierarchical topological order that constrains and organizes interface design.

Abstracting Global Interfaces To improve scalability and maintainability, recurring input–output patterns across modules are abstracted into common data structures or base classes, providing design anchors that enforce interface consistency and reduce redundancy. For example, algorithms can be unified under a `BaseEstimator` class to ensure standardized interaction with preprocessing and evaluation modules.

Adaptive Interface Design Within each file-level subgraph, leaf features are clustered into executable interfaces according to semantic relatedness. Independent features are implemented as standalone functions, while interdependent features are consolidated into shared classes with individual methods. For example, in Figure 2, `load_json` and `load_csv` are grouped into a `DataLoader` class, while regression variants are unified under `LinearModels`. This adaptive mapping balances granularity with cohesion, yielding a complete Repository Planning Graph (RPG) that preserves modularity and semantic consistency at the repository scale.

4 Graph-Guided Code Generation

As shown in Figure 1(D), given a user query and the completed RPG, ZeroRepo generates repositories by traversing the graph in topological order, ensuring dependencies are implemented before dependents. At each leaf node, test-driven development (TDD) is applied: a test is derived from the task specification, after which the corresponding functions or classes are implemented and validated against it; failing cases trigger revisions until the test passes. Only functions that pass all tests are committed to the repository, enabling incremental expansion while preserving stability. Additional details are provided in Appendix C.

Graph-Guided Localization and Editing To handle implementation and debugging requests, we adopt a two-stage workflow: first localizing the target in the RPG, then editing the associated code. Localization leverages the graph’s global structure and three complementary tools: (1) **RPG-Guided search**, which uses functionality-based fuzzy matching to identify candidate definitions; (2) **repository code view**, retrieving full interface bodies for inspection or modification; and (3) **dependency exploration**, tracing edges to reveal related modules and interactions. Once localized, the agent revises or generates the corresponding code to complete the requested implementation or repair.

Graph-Guided Test Validation To ensure correctness and contain errors early, validation follows a staged workflow aligned with the graph. Each function or class is first verified in isolation through unit tests automatically derived from its docstring. Validated components trigger regression tests upon modification, while completed subgraphs undergo integration tests to ensure consistent data flows and contracts across modules. A lightweight majority-vote diagnosis distinguishes genuine implementation errors from environment or test issues, automatically handling the latter and returning the former for repair through the localization–editing workflow.

5 Experiment Setup

Table 1: Overview of the six reference repositories and their paraphrased counterparts (Para. Name) in RepoCraft. #F. Cate. denotes functional categories, #Files the total source files, LOC the effective lines of code, and Task Counts the evaluation tasks for measuring code accuracy.

Real Repo	Para. Name	#F. Cate.	#Files	LOC	Code Tokens	Task Counts
scikit-learn	MLKit-Py	47	185	65,972	592,187	236
pandas	TableKit	81	217	106,447	943,873	175
sympy	SymbolicMath	40	699	218,924	2,339,881	192
statsmodels	StatModeler	88	271	83,325	893,824	234
requests	HttpEasy	22	17	2,793	22,297	50
django	PyWebEngine	42	681	109,457	917,622	165

5.1 RepoCraft Benchmark

A key challenge in evaluating repository-level generation is the absence of benchmarks that assess end-to-end reasoning and planning from scratch. Existing work either focuses on incremental development (editing, refactoring, or bug

fixing in existing codebases [26, 27, 35–37]) or repo generation but provides detailed skeletons and specifications that reduce the need for autonomous planning [38, 39]. RepoCraft addresses this gap by requiring agents to build complete repositories from high-level natural language descriptions and evaluating them against real-world projects in terms of scale, functionality, and correctness, with final statistics shown in Table 1.

5.1.1 Reference Repository Selection

To provide a strong reference for evaluating reasoning and planning, RepoCraft grounds assessment in six widely used Python projects: *scikit-learn*, *pandas*, *sympy*, *statsmodels*, *requests*, and *django*. These repositories exemplify high-quality engineering practice as they have been developed by active communities, exhibit modular structures, and include comprehensive test suites. Their diversity across scientific computing, data analysis, symbolic reasoning, web services, and full-stack frameworks ensures that the benchmark captures breadth and realism. To mitigate pretraining leakage, we paraphrase their names and descriptions into lexically distinct forms before providing them to agents.

5.1.2 Metrics

RepoCraft evaluates generated repositories along four complementary dimensions, with detailed definitions and formulas provided in Appendix D.3.1:

Functionality Coverage Coverage is measured as the proportion of functional categories, drawn from official documentation, that are represented by the generated functionalities. A category is counted as covered if at least one generated functionality corresponds to it. This metric reflects only the breadth of functionality, without assessing correctness. Reference taxonomies are provided in Appendix D.4.

Functionality Novelty Novelty is defined as the proportion of generated functionalities that fall outside the reference taxonomy, i.e., those assigned to the new features category. It captures the system’s ability to propose coherent but unseen capabilities beyond the ground specification.

Functionality Accuracy Accuracy evaluates correctness at the task level using two metrics: (1) **Pass Rate**, the fraction of ground-truth tests passed; and (2) **Voting Rate**, the fraction validated by majority-vote semantic checks. Unlike coverage, accuracy measures whether implementations faithfully realize the intended algorithms.

Code-Level Statistics We report repository scale indicators, including file count, normalized Lines of Code (LOC), and token count, measured after excluding non-core code such as tests and examples.

5.1.3 Functional Task Construction and Evaluation

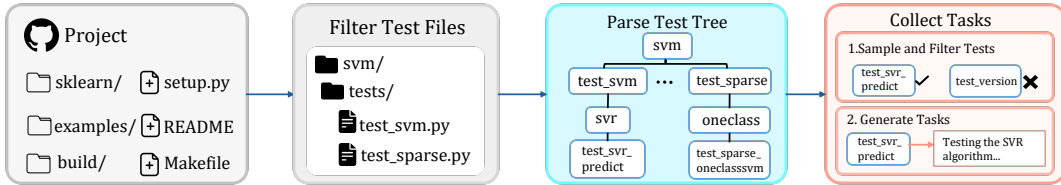


Figure 3: Pipeline for Evaluation Task Collection. It comprises test file filtering, hierarchical parsing into test trees, sampling and filtering, and final task generation.

To assess models’ planning ability on constructed repositories, we evaluate whether they (i) implement the intended algorithms and (ii) realize them correctly. Simple measures of repository size or coverage are insufficient for this purpose, so RepoCraft introduces task-level evaluations that capture both functional fidelity and implementation accuracy (see Appendix D.2 for details).

To enable such fine-grained evaluation, RepoCraft derives tasks from reference repositories. As shown in Figure 3, we collect all available test functions and classes, organize them hierarchically following each project’s modular structure, and apply stratified sampling to ensure representative coverage. Trivial or non-algorithmic tests are filtered out, resulting in a diverse and computationally meaningful set of 1,052 tasks that closely mirror practical software evaluation.

Each task includes a natural language description of the target algorithm, a ground-truth test, and auxiliary materials. Evaluation proceeds in three steps: (1) **Localization**, mapping requirements to candidate functions or classes in the generated repository; (2) **Semantic Validation**, applying majority-vote checks over two rounds to confirm fidelity to the specification; and (3) **Execution Testing**, adapting and running the ground-truth test to verify interface correctness under realistic inputs and outputs. This design mirrors real-world development while reducing sensitivity to spurious model errors. We use o3-mini as the base model for automated evaluation.

Table 2: Performance of agent frameworks and model backbones on RepoCraft. "Nov." denotes the novelty rate; the number in parentheses is Novel/Total, where Novel is the number of novel functionalities and Total is the total number of planned functionalities. Gold Projects are used as a confidence ablation for the automatic evaluation pipeline, and per-repository detailed results are reported in Appendix E.2.

Agent	Model	Cov. (%) \uparrow	Nov. (%) (Novel/Total) \uparrow	Pass. / Tot. (%) \uparrow	Files \uparrow	LOC \uparrow	Tokens \uparrow
MetaGPT	o3-mini	16.6	0.0 (0.0/24.8)	4.5 / 10.2	2.3	225.3	2180.3
	Qwen3-Coder	17.1	0.0 (0.0/32.7)	3.2 / 9.4	8.5	326.5	3369.3
ChatDev	o3-mini	18.3	9.2 (3.0/32.8)	2.6 / 10.5	5.8	410.3	4458
	Qwen3-Coder	22.1	3.9 (1.5/38.3)	6.9 / 11.6	6.3	540.7	5422.2
OpenHands	o3-mini	22.0	0.3 (0.1/36.5)	5.1 / 16.9	9.8	292.2	2712.8
	Qwen3-Coder	21.7	0.0 (0.0/33.7)	5.8 / 11.2	8.3	458.0	4778.3
Paper2Code	o3-mini	21.7	5.2 (2.1/40.0)	6.0 / 15.8	7.2	547.7	5920.8
	Qwen3-Coder	30.2	5.5 (4.0/73.8)	4.9 / 15.9	8.8	1365.2	14,555.0
Codex CLI	o3 pro	28.4	0.0 (0.0/48.5)	11.0 / 20.0	5.3	611.5	6248.5
Gemini CLI	gemini 2.5 pro	42.0	0.6 (0.8/132.7)	14.5 / 37.9	15.2	1484.8	14,922.2
Claude Code CLI	claude 4 sonnet	54.2	6.7 (41.6/621.0)	33.9 / 52.5	33.3	10,586.7	105,236.2
Gold Projects	Human Developers	-	-	81.0 / 92.0	345	97,819.7	951,614
ZeroRepo	o3-mini	81.5	13.6 (151.5/1114.2)	69.7 / 75.0	271.5	23,977.3	260,761.2
	Qwen3-Coder	75.1	9.2 (108.3/1173.3)	57.3 / 68.0	389.0	36,941.0	445,511.8

5.2 Baselines

We compare three paradigms: (1) **Multi-agent** frameworks (MetaGPT [7], ChatDev [8]) assigning specialized roles for end-to-end development; (2) **Workflow-based** system (Paper2Code [9]) with a fixed three-stage pipeline; (3) **Terminal agents** (Codex CLI [40], Claude Code CLI [12], Gemini CLI [13], OpenHands [11]) performing natural language editing, debugging, and multi-file reasoning. For comparability, MetaGPT, ChatDev, and Paper2Code are run with two backbones: o3-mini [41] and Qwen3-Coder-480B-A35B-Instruct (Qwen3-Coder) [42]. Codex CLI, Claude Code CLI, and Gemini CLI are evaluated with their official strongest model. **We enable Terminal Agents to retrieve real-world knowledge via web search.** To ensure fairness, all runs extend to 30 iterations, with agents prompted at each step to propose or implement functionality.

5.3 Implementation Details

We conduct 30 iterations for feature selection in Proposal-Level Graph Construction. Each function in the Code Generation Stage undergoes up to 8 debugging iterations, with 20 localization attempts per iteration. For test failures, we use 5-round majority voting to attribute and allow up to 20 remediation attempts for test or environment errors.

6 Main Results

RPG enables richer functionality and larger repositories. ZeroRepo demonstrates that RPG-guided planning yields repositories of substantially greater scale, diversity, and novelty than existing approaches. On RepoCraft, it achieves up to 81.5% coverage with o3-mini, representing a 27.3% absolute improvement over the strongest baseline (Claude Code at 54.2%). Beyond covering the required functionality, ZeroRepo also exhibits strong innovation, attaining novelty rates of 11–13% with over 100 new functionalities, whereas most baselines contribute fewer than 10. In terms of repository size, ZeroRepo with Qwen3-Coder generates 36K LOC and 445K tokens, corresponding to 3.9 \times the code scale of Claude Code and about 68 \times that of other baselines. Among these approaches, ZeroRepo is the closest to human-developed Gold Projects, underscoring that RPG serves as the key structured representation for building repositories that are larger, more diverse, and closer to real-world software complexity.

RPG enhances reasoning consistency and structural fidelity. Beyond scale, ZeroRepo delivers substantially higher correctness and stability. To ensure reliability, we first validate the automatic localization and validation pipeline on human-developed Gold Projects, where it achieves 81.0% pass rate and 92.0% voting agreement, establishing the ceiling under our test harness. Under the same protocol, ZeroRepo attains a 69.7% pass rate with o3-mini, an absolute improvement of 35.8% compared to the strongest baseline (Claude Code at 33.9%). These results indicate that RPG serves as a structured reasoning representation that enforces modular boundaries and functional contracts, thereby supporting coherent planning and yielding repositories that more faithfully realize intended specifications.

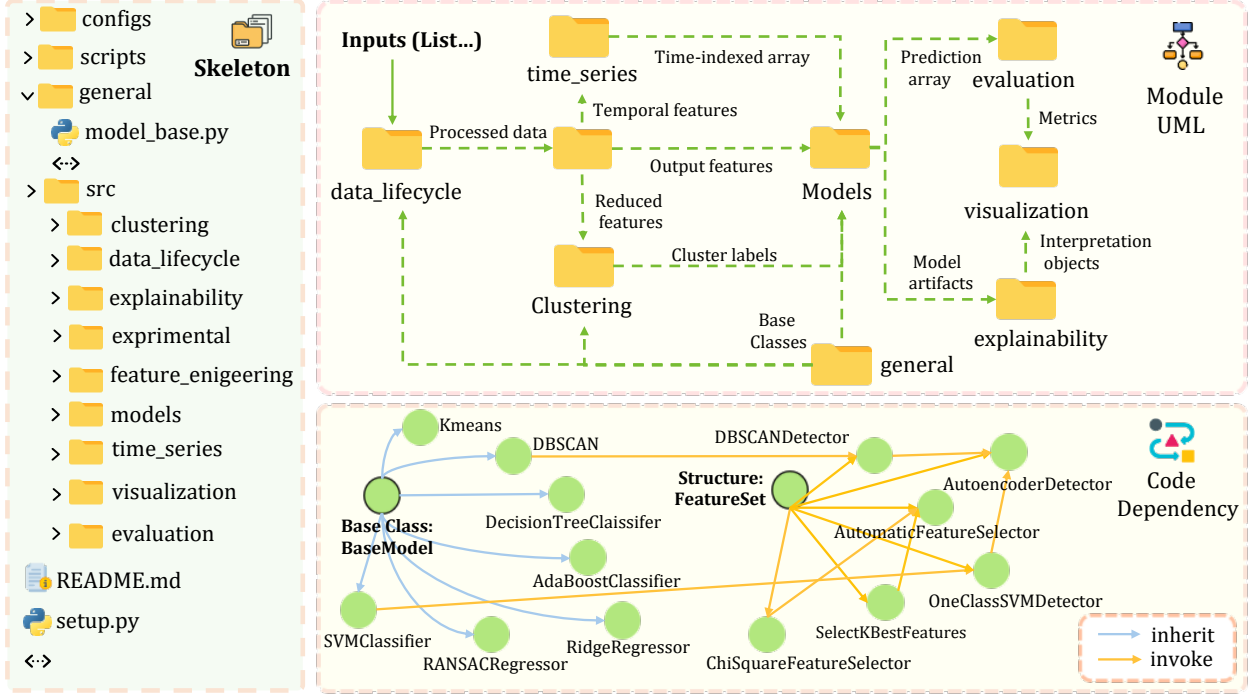


Figure 4: Illustration of dependencies in the repository generated by Qwen3-Coder on MLKit-Py, showing (1) the repository skeleton at the folder/file level, (2) inter-module data flows, and (3) class and function dependencies.

RPG induces complex data flows and dependencies. To illustrate the capacity of RPG-guided planning for generating complex repositories, we visualize ZeroRepo with Qwen3-Coder on the MLKit-Py task. At the file level, RPG organizes a coherent folder hierarchy; at the module level, inter-module flows define execution pipelines from `data_lifecycle` through `clustering` and `models` to `evaluation`; and at the function level, inheritance and invocation edges capture intricate class interactions. These results show that RPG induces layered dependencies and coordinated execution, enabling repositories with both structural complexity and internal coherence.

7 Analysis

7.1 Analysis of the RPG’s Scalability

RPG enables near-linear growth of repository functionalities.

A key question in repository-level generation is whether functionalities can continue to expand with iterative planning, or whether growth quickly stagnates. To examine this, we compute the number of planned features at each iteration on RepoCraft, averaging across 30 rounds for strong baselines (Claude Code, Gemini CLI, Codex CLI) and for ZeroRepo. As shown in Figure 5, ZeroRepo exhibits near-linear growth, surpassing 1,100 leaf features with o3-mini, while natural-language-based baselines display limited scalability: Claude Code grows steadily but with diminishing rates, Gemini CLI increases only slowly before converging by round 30, and Codex ceases to add features after just 4–5 iterations. These results demonstrate that the RPG provides a persistent and extensible planning substrate, enabling high-level goals to be refined into progressively richer functionalities. In contrast to natural-language representations, which degrade in coherence and stagnate, RPG sustains structural consistency and extensibility, establishing it as a stronger representational foundation for modeling repositories with increasingly complex functionalities and architectures.

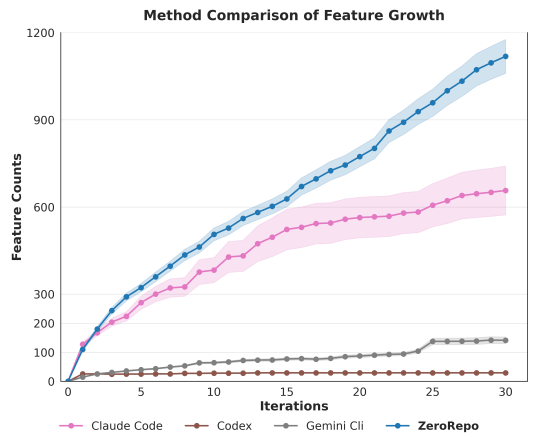


Figure 5: Feature comparison of ZeroRepo (o3-mini) against strong baselines (Codex, Gemini CLI, Claude Code) across iterations.

RPG ensures near-linear growth in repository size. Functional scalability is only meaningful if it translates into executable code. To examine this, we analyze how repository size evolves across iterations, measured in lines of code (LOC). As shown in Figure 6, ZeroRepo sustains near-linear growth, surpassing 30K LOC within 30 iterations. In contrast, natural-language-based baselines stagnate: Claude Code and Gemini CLI plateau around 3–4K LOC, while Codex remains below 1K LOC. This divergence reflects a fundamental representational gap. Natural language planning rapidly accumulates inconsistencies across iterations, leading to fragmented specifications that fail to translate into coherent code. In contrast, the RPG maintains a persistent, extensible structure in which proposed functionalities are grounded in explicit modules, interfaces, and data flows. This grounding ensures that planned expansions are consistently realized as executable code, producing repositories that grow not only in size but also in organizational coherence. These results highlight the RPG’s ability to sustain repository scaling in both volume and integrity, positioning it as a robust representational basis for long-horizon code generation.

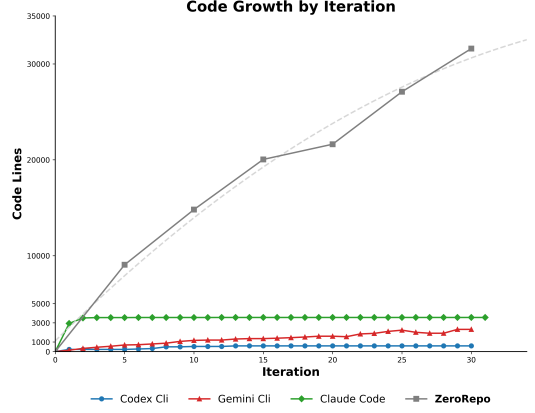


Figure 6: Scaling behavior of total lines of code across iteration steps on MLKit-Py

7.2 Analysis of RPG’s Stability and Innovation Potential

RPG supports comprehensive and extensible functionality. A central challenge in repository-level generation is ensuring repositories not only satisfy user-specified requirements but also extend beyond them coherently. As shown in Table 3, ZeroRepo steadily increases coverage from 70.2% at 5 iterations to nearly 96% at 30, far surpassing baselines below 60% (Table 2). Simultaneously, it maintains novelty, reaching 8% with over 100 additional features, whereas baselines contribute fewer than 50. These results suggest that RPG functions as a persistent planning substrate, enabling repositories to achieve comprehensive coverage while supporting principled growth beyond reference implementations. Representative examples in Appendix E.3 validate that RPG sustains coherence in both coverage allocation and novel feature proposals.

Table 3: Coverage and Novelty of the Constructed RPG over Iterations on MLKit-Py (o3-mini)

Iteration	Cov. (%) ↑	Nov. (%) ↑
5	70.2	4.6 (15.3/336.1)
10	80.9	5.4 (29.01/542.0)
15	83.0	4.9 (39.0/796.0)
20	85.1	5.2 (51.0/981.0)
25	87.2	7.0 (73.5/1043.0)
30 (ours)	95.7	7.9 (99.4/1258.0)

7.3 Analysis of Graph-Guided Localization

Graph guidance accelerates agent localization. We evaluate the impact of RPG guidance by comparing localization steps with and without graph support (Table 4). Across Integration Testing (Int. Test.), Incremental Development (Incr. Dev.), and Debugging (Debug.), graph-guided search reduces effort by 30–50%. This demonstrates that RPG equips agents with a principled navigation mechanism, enabling faster dependency tracing, more accurate bug localization, and smoother module integration, thereby improving repository development efficiency. Compared to natural language, RPG offers a global structural representation of the repository, enabling agents to localize targets from a functionality-wide perspective and accelerating the development cycle.

Table 4: Ablation of Graph-Guided Localization on MLKit-Py (o3-mini). Steps (mean ± SD). “wo/Graph” denotes ZeroRepo without Graph.

Category	IntTest	IncDev	Debug
ZeroRepo	6.2 ± 2.1	6.8 ± 1.8	5.8 ± 2.8
- w/o Graph	13.3 ± 11.1	10.8 ± 2.6	8.5 ± 2.9

8 Conclusion

In this paper, we presented the Repository Planning Graph (RPG), a structured representation that unifies proposal- and implementation-level planning for repository generation. Built on RPG, we introduced ZeroRepo, a graph-driven framework that achieves state-of-the-art coverage, correctness, and scalability on the RepoCraft benchmark. Our analyses show that RPG models complex dependencies, enables progressively more sophisticated planning through near-linear scaling of functionality and code size, and improves agents’ repository understanding, thereby accelerating localization. These results highlight the potential of graph-based representations as a foundation for advancing long-horizon and large-scale repository generation.

References

- [1] Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv preprint arXiv:2406.11931*, 2024.
- [2] Yaoliang Wang, Haoling Li, Xin Zhang, Jie Wu, Xiao Liu, Wenxiang Hu, Zhongxin Guo, Yangyu Huang, Ying Xin, Yujiu Yang, et al. Epicoder: Encompassing diversity and complexity in code generation. *arXiv preprint arXiv:2501.04694*, 2025.
- [3] Yifei Liu, Li Lyna Zhang, Yi Zhu, Bingcheng Dong, Xudong Zhou, Ning Shang, Fan Yang, and Mao Yang. rstar-coder: Scaling competitive code reasoning with a large-scale verified dataset. *arXiv preprint arXiv:2505.21297*, 2025.
- [4] Huaye Zeng, Dongfu Jiang, Haozhe Wang, Ping Nie, Xiaotong Chen, and Wenhui Chen. Acecoder: Acing coder rl via automated test-case synthesis. *arXiv preprint arXiv:2502.01718*, 2025.
- [5] Hongyuan Tao, Ying Zhang, Zhenhao Tang, Hongen Peng, Xukun Zhu, Bingchang Liu, Yingguang Yang, Ziyin Zhang, Zhaogui Xu, Haipeng Zhang, et al. Code graph model (cgm): A graph-integrated large language model for repository-level software engineering tasks. *arXiv preprint arXiv:2505.16901*, 2025.
- [6] Haiyang Li. Mrg-bench: Evaluating and exploring the requirements of context for repository-level code generation. *arXiv preprint arXiv:2508.02998*, 2025.
- [7] Sirui Hong, Xiwu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 3(4):6, 2023.
- [8] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023.
- [9] Minju Seo, Jinheon Baek, Seongyun Lee, and Sung Ju Hwang. Paper2code: Automating code generation from scientific papers in machine learning. *arXiv preprint arXiv:2504.17192*, 2025.
- [10] Zijie Lin, Yiqing Shen, Qilin Cai, He Sun, Jinrui Zhou, and Mingjun Xiao. Autop2c: An llm-based agent framework for code repository generation from multimodal content in academic papers. *arXiv preprint arXiv:2504.20115*, 2025.
- [11] Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024.
- [12] Anthropic. Claude code: agentic coding command-line tool, 2025. URL <https://www.anthropic.com/claude-code>.
- [13] Google. Gemini cli: Open-source ai agent for developer terminals, June 2025. Accessed: 2025-07-22.
- [14] Evan Wang, Federico Cassano, Catherine Wu, Yunfeng Bai, Will Song, Vaskar Nath, Ziwen Han, Sean Hendryx, Summer Yue, and Hugh Zhang. Planning in natural language improves llm search for code generation. *arXiv preprint arXiv:2409.03733*, 2024.
- [15] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 17682–17690, 2024.
- [16] Haotian Sun, Yuchen Zhuang, Linghai Kong, Bo Dai, and Chao Zhang. Adapllanner: Adaptive planning from feedback with language models. *Advances in neural information processing systems*, 36:58202–58245, 2023.
- [17] Yueheng Zhu, Chao Liu, Xuan He, Xiaoxue Ren, Zhongxin Liu, Ruwei Pan, and Hongyu Zhang. Adacoder: An adaptive planning and multi-agent framework for function-level code generation. *arXiv preprint arXiv:2504.04220*, 2025.
- [18] Amr Almorsi, Mohammed Ahmed, and Walid Gomaa. Guided code generation with llms: A multi-agent framework for complex code tasks. In *2024 12th International Japan-Africa Conference on Electronics, Communications, and Computations (JAC-ECC)*, pages 215–218. IEEE, 2024.
- [19] Nazmus Ashrafi, Salah Bouktif, and Mohammed Mediani. Enhancing llm code generation: A systematic evaluation of multi-agent collaboration and runtime debugging for improved accuracy, reliability, and latency. *arXiv preprint arXiv:2505.02133*, 2025.

- [20] OpenAI. Gpt-4o system card. arXiv preprint arXiv:2410.21276, 2024. URL <https://arXiv.org/abs/2410.21276>.
- [21] Anthropic. Introducing claude 4, May 2025. URL <https://www.anthropic.com/news/claude-4>.
- [22] Gemini Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, June 2025. URL https://storage.googleapis.com/deepmind-media/gemini/gemini_v2_5_report.pdf.
- [23] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [24] Muhammad Nouman Zafar, Wasif Afzal, and Eduard Enoiu. Evaluating system-level test generation for industrial software: A comparison between manual, combinatorial and model-based testing. In *Proceedings of the 3rd ACM/IEEE international conference on automation of software test*, pages 148–159, 2022.
- [25] Arghavan Moradi Dakhel, Amin Nikanjam, Vahid Majdinasab, Foutse Khomh, and Michel C Desmarais. Effective test generation using pre-trained large language models and mutation testing. *Information and Software Technology*, 171:107468, 2024.
- [26] Dhruv Gautam, Spandan Garg, Jinu Jang, Neel Sundaresan, and Roshanak Zilouchian Moghaddam. Refactorbench: Evaluating stateful reasoning in language agents through code. *arXiv preprint arXiv:2503.07832*, 2025.
- [27] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- [28] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
- [29] Jianwen Luo, Yiming Huang, Jinxiang Meng, Fangyu Lei, Shizhu He, Xiao Liu, Shanshan Jiang, Bin Dong, Jun Zhao, and Kang Liu. Gate: Graph-based adaptive tool evolution across diverse tasks. *arXiv preprint arXiv:2502.14848*, 2025.
- [30] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.
- [31] OpenAI. Introducing codex: a cloud-based software engineering agent, 2025. URL <https://openai.com/index/introducing-codex/>.
- [32] Google. Gemini cli: your open-source ai agent, 2025. URL <https://blog.google/technology/developers/introducing-gemini-cli-open-source-ai-agent/>.
- [33] Karthik Valmeekam, Sarath Sreedharan, Matthew Marquez, Alberto Olmo, and Subbarao Kambhampati. On the planning abilities of large language models (a critical investigation with a proposed benchmark). *arXiv preprint arXiv:2302.06706*, 2023.
- [34] Ma’ayan Armony, Albert Meroño-Peñuela, and Gerard Canal. How far are llms from symbolic planners? an nlp-based perspective. *arXiv preprint arXiv:2508.01300*, 2025.
- [35] Linghao Zhang, Shilin He, Chaoyun Zhang, Yu Kang, Bowen Li, Chengxing Xie, Junhao Wang, Maoquan Wang, Yufan Huang, Shengyu Fu, et al. Swe-bench goes live! *arXiv preprint arXiv:2505.23419*, 2025.
- [36] Wei Li, Xin Zhang, Zhongxin Guo, Shaoguang Mao, Wen Luo, Guangyue Peng, Yangyu Huang, Houfeng Wang, and Scarlett Li. Fea-bench: A benchmark for evaluating repository-level code generation for feature implementation. *arXiv preprint arXiv:2503.06680*, 2025.
- [37] Yiming Huang, Jianwen Luo, Yan Yu, Yitong Zhang, Fangyu Lei, Yifan Wei, Shizhu He, Lifu Huang, Xiao Liu, Jun Zhao, et al. Da-code: Agent data science code generation benchmark for large language models. *arXiv preprint arXiv:2410.07331*, 2024.
- [38] Wenting Zhao, Nan Jiang, Celine Lee, Justin T Chiu, Claire Cardie, Matthias Gallé, and Alexander M Rush. Commit0: Library generation from scratch. *arXiv preprint arXiv:2412.01769*, 2024.
- [39] Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin, Rachel Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, et al. Paperbench: Evaluating ai’s ability to replicate ai research. *arXiv preprint arXiv:2504.01848*, 2025.
- [40] OpenAI. Codex cli: Local-first terminal ai coding agent. <https://openai.com/introducing-codex>, April 2025. Launched April 15, 2025 as an open-source CLI coding agent that executes commands locally :contentReference[oaicite:1]index=1.

- [41] OpenAI. Openai o3-mini reasoning model. System Card and official release (OpenAI), 2025. URL <https://openai.com/index/openai-o3-mini/>. Released January 31, 2025. Accessed: 2025-07-28.
- [42] Qwen Team. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.

A Appendix of Proposal-Level Graph Construction

The construction of the RPG is central to our framework, as it transforms high-level repository goals into a structured and persistent representation. The process starts with carefully designed prompting strategies for selecting repository-relevant functionalities from the global feature ontology, followed by iterative refinement to ensure both semantic coverage and modular coherence.

A.1 Algorithms of Functionality Graph Construction

Algorithm 1 Feature Sampling with Diversity-Aware Rejection Sampling

Require: Root node R ; frequency library F ; temperature t ; per-tree sample size S ; overlap threshold ρ ; maximum number of retries T_{\max}

```

1: function BASESAMPLE( $R, F, t, S$ )
2:   selected_set  $\leftarrow \emptyset$ 
3:   for  $s = 1$  to  $S$  do
4:      $C \leftarrow \text{GET\_CHILDREN}(R)$ 
5:     if  $C = \emptyset$  then
6:       break
7:     end if
8:      $f_i \leftarrow F[i]$  for all  $i \in C$ 
9:      $p_i \leftarrow f_i / \sum_{j \in C} f_j$  for all  $i \in C$ 
10:     $q_i \leftarrow \text{TEMPTRANSFORM}(p_i, t)$  for all  $i \in C$ 
11:    cur_node  $\leftarrow \text{SAMPLE\_NODE}(C, [q_1, q_2, \dots])$ 
12:    selected_set.ADD(cur_node)
13:     $R \leftarrow \text{cur\_node}$  ▷ move root downward for next step
14:   end for
15:   return selected_set
16: end function

17: function REJECTSAMPLE( $R, F, t, S, \rho, T_{\max}$ )
18:   best_T  $\leftarrow \emptyset$ ; best_ovl  $\leftarrow +\infty$ 
19:    $T^* \leftarrow \emptyset$ 
20:   for  $\tau = 1$  to  $T_{\max}$  do ▷ retry up to  $T_{\max}$  times
21:      $T_{\text{cand}} \leftarrow \text{BASESAMPLE}(R, F, t, S)$  ▷ sample a candidate tree
22:     ovl  $\leftarrow \text{OVERLAP}(T_{\text{cand}}, \mathcal{S}_{\text{seen}})$  ▷ compute overlap with seen nodes
23:     if ovl  $<$  best_ovl then
24:       best_ovl  $\leftarrow$  ovl; best_T  $\leftarrow T_{\text{cand}}$  ▷ update best candidate so far
25:     end if
26:     if ovl  $\leq \rho$  then
27:        $T^* \leftarrow T_{\text{cand}}$  ▷ accept immediately if overlap  $\leq$  threshold
28:       break
29:     end if
30:   end for
31:   if  $T^* = \emptyset$  then
32:      $T^* \leftarrow \text{best\_T}$  ▷ fallback: choose least-overlap candidate
33:   end if
34:   return  $T^*$  ▷ return the final selected tree
35: end function

```

Rejection Sampling Algorithm We extend the base sampling strategy introduced in EPICODER [2] by incorporating a diversity-aware rejection mechanism, as shown in Algorithm 1. At each step, a candidate tree is accepted only if its overlap with previously sampled nodes is below a specified threshold; otherwise, the tree with the minimal overlap is returned. This encourages broader feature space exploration under a limited number of sampling iterations.

Repository-Aligned Subtree Selection Algorithm 2 outlines the procedure for constructing a repository-specific feature subtree from a global feature tree. The algorithm iteratively selects candidate features based on a combination of exploitation (retrieving top-scoring nodes) and exploration (sampling unvisited nodes). At each iteration, an LLM agent

filters and ranks candidates, proposes missing but relevant features, and performs batch-level self-checks to ensure internal consistency. Accepted candidates are incorporated into the current subtree, and the process continues until a fixed iteration budget is reached. The resulting subtree captures features most relevant to the target repository while balancing coverage and quality.

Algorithm 2 Repository-Specific Subtree Selection

Require: Global Feature Tree \mathcal{T} ; Repo description \mathcal{R} ; iteration budget K ; top- k ; termination threshold τ ; batch size B ; LLM

Ensure: Repository-specific subtree \mathcal{T}'

- 1: Initialize current repo tree $\mathcal{T}' \leftarrow \emptyset$; missing features $\mathcal{C}_{\text{missing}} \leftarrow \emptyset$; visited set $\mathcal{V} \leftarrow \emptyset$
- 2: **for** $k = 1 \dots K$ **do** ▷ iterate with given budget
- 3: $\mathcal{E}_{\text{exploit}} \leftarrow \text{RETRIEVETOPK}(\mathcal{T}, \mathcal{R}, k = \text{top-}k)$ ▷ select promising nodes (exploit)
- 4: $\mathcal{E}_{\text{explore}} \leftarrow \text{SAMPLEUNVISITED}(\mathcal{T}, \mathcal{V})$ ▷ sample unexplored nodes (explore)
- 5: **// Candidate selection via LLM**
- 6: $\mathcal{C}_{\text{exploit}} \leftarrow \text{LLM.SELECTEXPLOITCANDIDATES}(\mathcal{E}_{\text{exploit}}, \mathcal{T}', \mathcal{R})$ ▷ filter exploit candidates
- 7: $\mathcal{C}_{\text{explore}} \leftarrow \text{LLM.SELECTEXPLOREcandidates}(\mathcal{E}_{\text{explore}}, \mathcal{T}', \mathcal{R})$ ▷ filter explore candidates
- 8: $\mathcal{C}_{\text{missing}} \leftarrow \mathcal{C}_{\text{missing}} \cup \text{LLM.PROPOSEMISSING}(\mathcal{T}', \mathcal{R})$ ▷ generate novel candidates not in tree
- 9: $\mathcal{C}_{\text{raw}} \leftarrow \mathcal{C}_{\text{exploit}} \cup \mathcal{C}_{\text{explore}} \cup \mathcal{C}_{\text{missing}}$ ▷ merge all candidate sources
- 10: **// Batch self-check (filter useful paths within each batch)**
- 11: **for all** batch $\mathcal{B} \subseteq \mathcal{C}_{\text{raw}}$ with $|\mathcal{B}| \leq B$ **do** ▷ process in small batches
- 12: $\mathcal{B}^* \leftarrow \text{LLM.SELFCHECK}(\mathcal{T}', \mathcal{B})$ ▷ accept only consistent/relevant paths
- 13: $\mathcal{T}' \leftarrow \text{INSERTPATHS}(\mathcal{T}', \mathcal{B}^*)$ ▷ expand repo-specific tree
- 14: $\mathcal{V} \leftarrow \mathcal{V} \cup \mathcal{B}$ ▷ mark all evaluated paths as visited
- 15: **end for**
- 16: **end for**
- 17: **return** \mathcal{T}' ▷ return final subtree

Repository Subtree Reorganization into the functionality graph The algorithm operates in three stages to refactor subtree. In the first stage, an LLM agent iteratively extracts features from the input, organizing them into subgraphs until sufficient coverage of leaf nodes is reached. In the second stage, the agent reorganizes subgraphs by merging semantically related components or moving branches across groups to improve structure. Finally, each subgraph is refined to ensure naming consistency and hierarchical coherence, yielding a clean, interpretable functionality graph.

A.2 Detailed Construction Process

Global Feature Tree The global feature tree consists of more than one million nodes across seven hierarchical levels (Table 5), reflecting a broad and diverse functional knowledge base. Nevertheless, the distribution of features across Level-1 categories is highly skewed (Figure 8). In particular, the *data processing* branch dominates the tree, while many other categories contain only a small number of nodes, resulting in a pronounced long-tail distribution. Such bias is inherent to real-world software ecosystems, where data processing utilities are disproportionately prevalent compared to specialized functionalities. As a consequence, constructing a repository-specific RPG requires large-scale filtering and reorganization in order to extract the most relevant features and mitigate the imbalance of the global distribution.

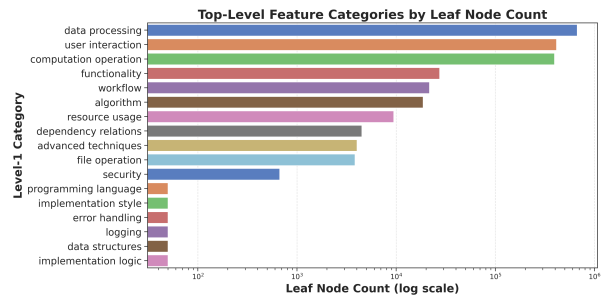


Figure 8: Distribution of feature counts under Level-1 categories in the global feature tree.

Model-Specific Growth Patterns Beyond the two traces in Fig. 7. Concretely, *qwen3-coder* exhibits the most open expansion, with a linear increase in leaf counts—maximizing coverage early but with a higher risk of admitting loosely related features. *o3-mini* follows with a moderately aggressive trajectory, striking a balance between breadth and relevance. Together, these curves delineate points on the recall-precision spectrum of subtree selection strategies matched to repository needs.

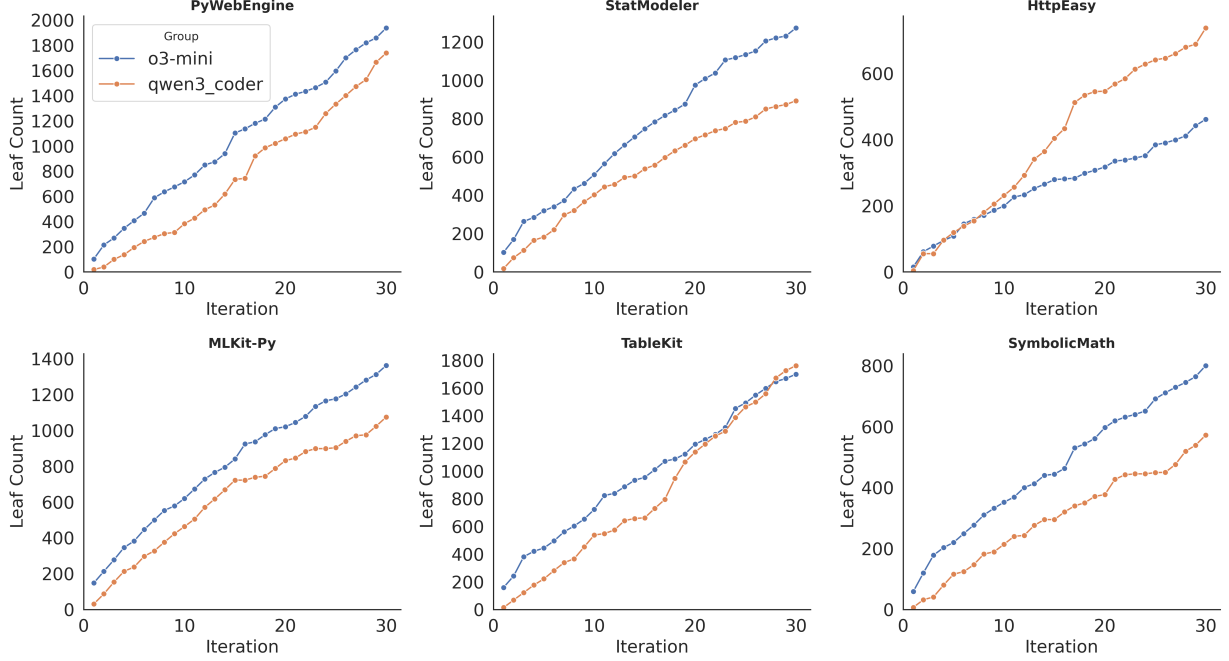


Figure 7: Evolution of Feature Tree Leaf Counts over Iterations Across Repositories, Highlighting the Differences Between qwen3 coder and o3-mini

Table 5: Statistics of the global feature tree across hierarchical levels, with representative examples from each level.

Level	#Elements	Examples
1	17	functionality, data structures, data processing
2	1,527	text processing, process monitoring, create flowchart
3	21,739	heap allocation, dayjs, affine transformation
4	113,348	update record by ID, automated versioning, M5 Model Tree
5	613,311	add vertices, angular velocity, find minimum by group, mark outlier data
6	33,801	min with inclusion, multiple with keyword
7	781	validate against thesaurus, swipe event detection

From Global to Repository-Specific Distributions The comparison between the global feature tree (Fig.8) and the final repository-specific profiles (Figs.10) highlights the transformative effect of model-guided reorganization. While the global tree is dominated by generic categories such as data processing and user interaction, the restructured graphs consistently downweight these high-frequency but less discriminative categories and elevate domain-relevant branches to the foreground. This shift effectively counteracts the inherent long-tail bias of the global ontology, redistributing feature density toward categories that better capture repository semantics. As a result, the constructed graphs are not only semantically sharper but also more functionally coherent with respect to the target domain. Between models, *qwen3-coder* favors broad coverage with slower convergence and higher variance, whereas *o3-mini* achieves a more balanced trade-off between generality and specificity. Together, these contrasting tendencies illustrate complementary strategies along the recall-precision spectrum, offering flexibility in matching feature selection to downstream repository needs.

Final Graph Structures The final RPGs (Figure 11a, 11b) reveal how repository-specific functionalities are consolidated into coherent modular organizations. Compared to the more diffuse subtree distributions, the resulting graphs exhibit a markedly skewed allocation of functionalities across subgraphs: a small number of core subgraphs absorb the majority of features, while peripheral subgraphs remain lightweight. This reflects a natural modularization process, where dominant clusters correspond to central repository capabilities and minor clusters capture auxiliary or specialized functions. Between models, the partitioning strategies diverge: *qwen3-coder* produces a larger number of medium-sized subgraphs, favoring breadth and parallel coverage; whereas *o3-mini* yields a more balanced distribution, with several

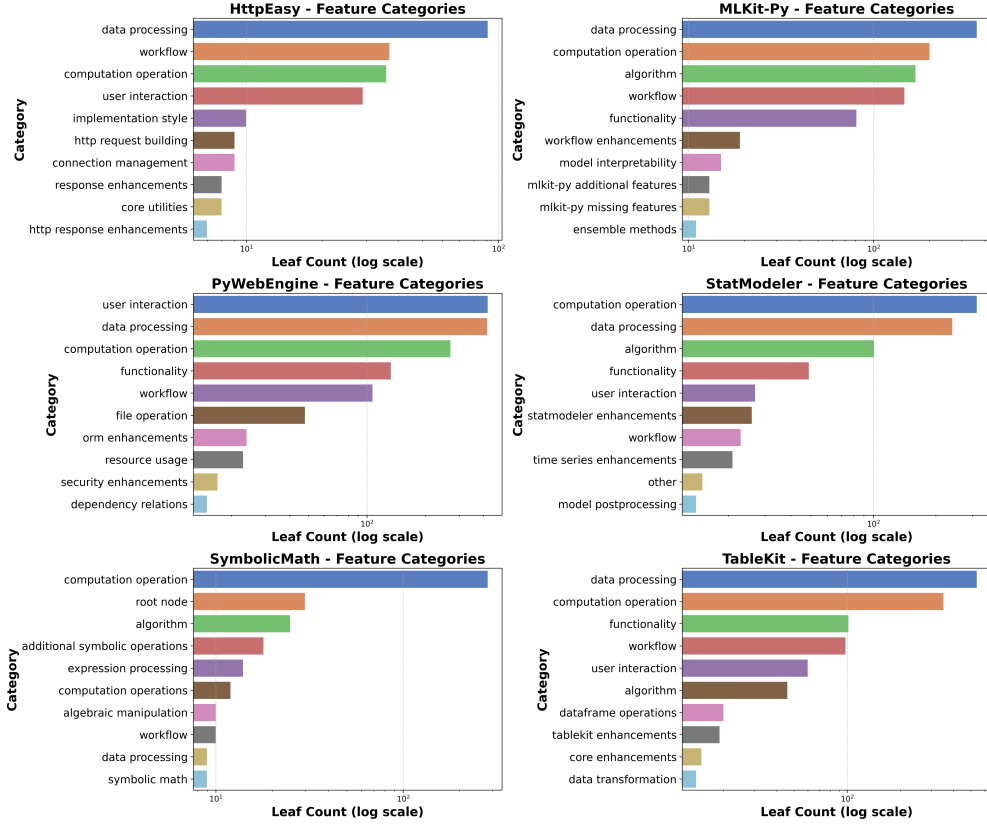


Figure 9: Final distribution of feature counts across subtrees for all repositories under *o3-mini*. The figure shows how features are reorganized after the iterative construction process, reflecting the model’s preference in balancing breadth and precision.

subgraphs of comparable size anchoring distinct semantic roles. These differences indicate that model-driven reorganization not only mitigates the global ontology’s long-tail bias but also shapes the granularity of modular decomposition, thereby influencing how functional responsibilities are distributed within the generated graph.

A.3 Prompt Template

Parts of Prompt Templates for Exploit–Explore Strategy in Subtree Selection

Prompt for Exploitation Paths

You are a GitHub project assistant responsible for expanding a repository’s feature tree through path-based modifications to ensure alignment with the project’s goals.

In each response, you will be given:

- An Exploit Feature Tree: A curated subset of high-relevance feature paths.
- The Current Repository Feature Tree.

When returning selected paths, always use “path/to/feature” format with ‘/’ as the separator.

Objective (excerpt)

- Expand the Repository Feature Tree so it: 1. Aligns with the repository’s purpose and scope.
- Achieves broad coverage across functional areas.
- Ensures essential capabilities are represented.
- Identify and fill critical gaps. ...

Selection Principles (excerpt)

- Select exclusively from the Exploit Feature Tree.
- Include all non-duplicated, useful paths.

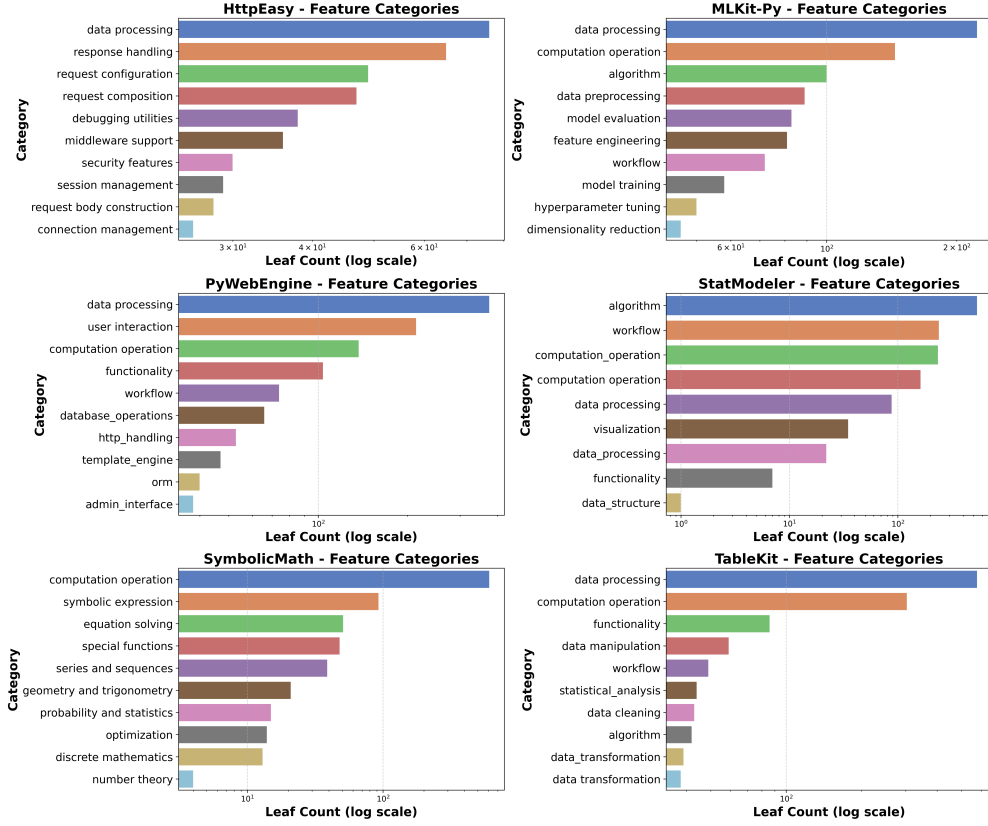


Figure 10: Final distribution of feature counts across subtrees for all repositories under *qwen3-coder*. The figure shows how features are reorganized after the iterative construction process, reflecting the model’s preference in balancing breadth and precision.

- Maintain structural balance by covering underrepresented modules.

...

Exclusions (excerpt)

Skip generic infra (e.g., logging, configuration) and abstract goals (e.g., “optimize CPU usage”).

Response Format

Respond only with a Thought and an Action.

<think>

Reason about relevance and gaps in the Exploit Tree. </think>

<action>

```
{
  "all_selected_feature_paths": [
    "path/to/feature", ...
  ]
}
```

</action>

Prompt for Exploration Paths

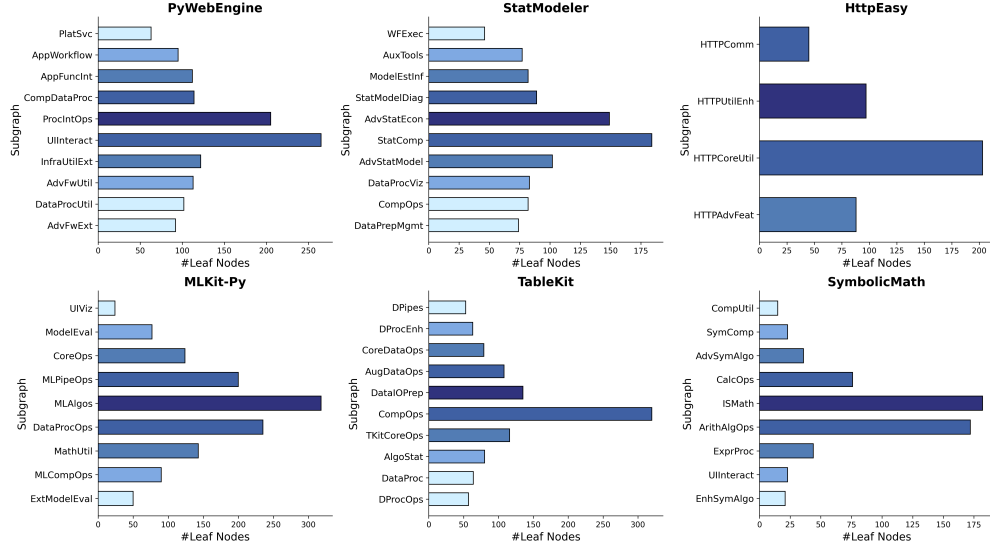
You are a GitHub project assistant responsible for expanding a repository’s feature tree through path-based modifications to ensure alignment with the project’s goals.

In each response, you will be given:

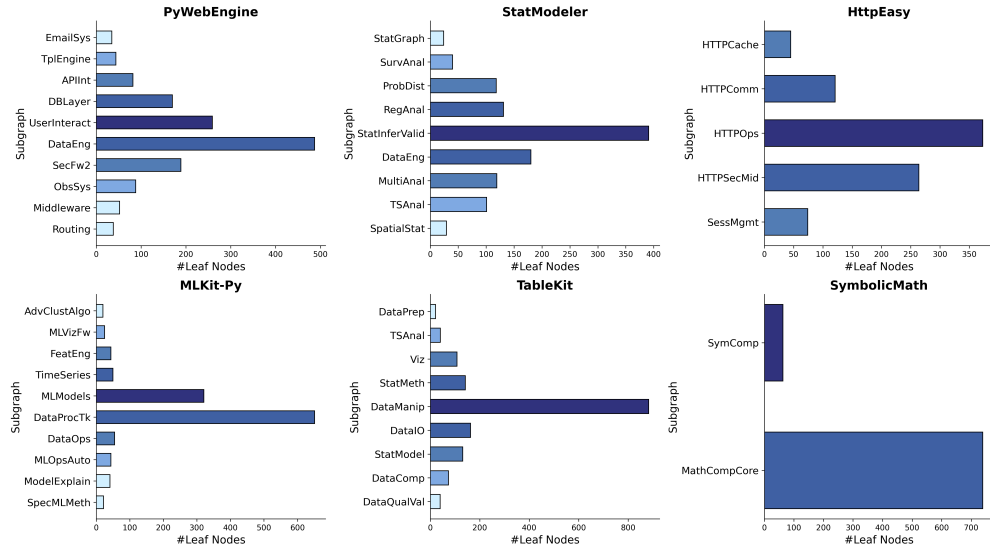
- A Sampled Feature Tree (Exploration Tree).
- The Current Repository Feature Tree.

When returning selected paths, always use “path/to/feature” format with ‘/’ as the separator.

Objective (excerpt)



(a) o3-mini



(b) qwen3-coder

Figure 11: Leaf node counts distribution across feature subgraphs in each repository RPG, reorganized by different models.

Improve and expand the Repository Feature Tree so that it: - Aligns with the repository's purpose and usage scenarios.

- Achieves comprehensive coverage of core and supporting areas.

...

Selection Principles (excerpt)

- Select only from the Exploration Tree.

- Include actionable, domain-relevant features.

- Skip paths already present in the current Repository Tree.

- Slight over-inclusion is acceptable.

...

Exclusions (excerpt)

Do not select generic infra (e.g., logging, config) or large-scale features (e.g., cloud integrations).

Response Format

Respond only with a single `<think>` and `<action>` block.

`<think>`

Explain how each Exploration Tree path was evaluated and why it was included or excluded.

`</think>`

`<action>`

```
{ "all_selected_feature_paths": [ "path/to/feature", ... ] }
```

`</action>`

Parts for Prompt Template for Retrieving Missing Features

Instruction You are a GitHub project assistant tasked with designing a functionally complete, production-grade repository.

Your goal is to identify and recommend **missing functional capabilities or algorithms** that the project should include, based on its real-world purpose, scope, and domain expectations.

Focus on intended functionality — not the existing Feature Tree, which may be incomplete.

Objective (excerpt)

Identify groups of functionally concrete features that: 1. Align with the repository’s domain and purpose.

2. Are entirely missing or only superficially represented.

3. Are specific and implementable (e.g., functions, classes, modules, algorithms).

Inclusion Criteria (excerpt)

- Must be code-level operations (computation, transformation, algorithm, evaluation).

- Realistically implementable within the repository’s scope.

- Both standard and advanced algorithms are allowed.

Exclusion Criteria (excerpt)

Do not include abstract intentions (e.g., “improve accuracy”), generic infra (e.g., logging, connectors), placeholders, or duplicates.

Naming Rules (excerpt)

- Use 3–5 lowercase words, separated by spaces.

- Each leaf node must describe a concrete algorithm or behavior.

- Avoid vague terms, camelCase, or snake_case.

Structure Guidelines (excerpt)

- Organize into logical hierarchies (up to 4–5 levels).

- Reflect computational architecture, not documentation taxonomy.

Response Format Respond with ONLY a `<think>` and `<action>` block:

`<think>`

Reason about functional domains, workflows, and algorithms that are missing from the current Feature Tree but expected in real-world use.

`</think>`

`<action>`

```
{ "missing_features": { "root node": { "child node 1": [ "leaf feature 1", "leaf feature 2" ], "child node 2": [ "leaf feature 3", "leaf feature 4" ] } } }
```

`</action>`

B Appendix of Implementation-Level Graph Construction

This section illustrates how the RPG is enriched with file organization and function design to form concrete code structures.

B.1 Prompt Template for Implementation-Level Graph Construction

We provide the prompt templates that guide the transformation from graph subtrees into modular code skeletons.

Prompt

You are a system architect tasked with designing the inter-subtree data flow for a Python software repository. Your goal is to define how data moves between functional modules (subtrees) — including who produces it, who consumes it, and how it is transformed — and express this as a structured, directed graph.

Data Flow

Format

```
[ { "from": "<source subtree name>", "to": "<target subtree name>", "data_id": "<unique name or description of the data>", "data_type": "<type or structure of the data>", "transformation": "<summary of what happens to the data, if anything>" }, ... ]
```

Validity & Structural Constraints

2. Full Connectivity Required

- Every subtree listed in {trees_names} must appear in at least one edge.
- No subtree should be isolated or unused.

3. Acyclic Structure

- The data flow must form a Directed Acyclic Graph (DAG):

4. Field Guidelines

- 'data_id': Use unique, descriptive names to identify each data exchange.
- 'data_type': Use precise and interpretable types to describe the structure, format, or abstract role of the data being passed.
- 'transformation': Describe how the data is modified, filtered, enriched, or combined. If unchanged, say "none".

...

Output Format

<solution>

```
[ { "from": "...", "to": "...", "data_id": "...", "data_type": "...", "transformation": "..." }, ... ] </solution>
```

Parts of Prompt Templates for Raw Skeleton Mapping

You are a repository architect responsible for designing the initial project structure of a software repository in its early development stage. Your task is to design a clean, modular file system skeleton that organizes the repository into appropriate top-level folders based on these subtrees.

Requirements

1. The folder structure must clearly separate each functional subtree and reflect logical domain boundaries.
2. Folder names must be concise, meaningful, and follow Python conventions (e.g., 'snake_case'). Names should feel natural and developer-friendly.
3. Folder names do not need to match subtree names exactly.
 - Treat subtree names as functional labels.
 - Rename folders as needed for clarity and convention, while preserving the correct mapping.
 - When assigning a subtree to a folder, include the exact subtree name in the mapping (e.g., "ml_models": ["Machine Learning"]).
4. You may choose a flat structure (all folders at root level) or a nested structure (e.g., under 'src'), depending on what best supports clarity, organization, and practical use.
5. Include commonly used auxiliary folders as appropriate.
6. The proposed structure should balance clarity, scalability, and maintainability. Avoid unnecessary complexity or excessive nesting.

...

Output Format

Return a single JSON-style nested object representing the repository folder structure:

- `"folder_name": ["Subtree Name"]` means this folder is assigned to a specific subtree. The name in the list must match exactly what appears in the given list of subtrees.
- `"folder_name": []` means the folder exists but does not correspond to a specific subtree (e.g., utility or support folders).
- `"file_name.ext": null` indicates the presence of a file. File content is not required.

Prompt for Mapping Feature Paths to Python Skeleton Files

You are a repository architect tasked with incrementally assigning all remaining leaf-level features from a functional subtree into the repository's file structure. This is an iterative process. You are not expected to assign all features at once — each round should make clear, meaningful progress. Your ultimate goal is a production-grade file organization that evolves cleanly and logically over time.

Context

In each iteration, you will receive:

- A list of unassigned leaf features (each is a full path like `"a/b/c"`).
- A designated functional folder under which all new paths must begin.
- A partial skeleton showing the current structure (existing assignments are hidden).

Assign the remaining features to `.py` file paths that:

- Begin with the designated folder.
- Group semantically related features together.
- Reflect how real developers would modularize logic in a production Python codebase. - Prefer organizing major functional categories into subfolders when appropriate.

File & Folder Structure

- Group features by functionality into logically meaningful modules that reflect real-world development practice.
- Avoid bundling many unrelated features into a single file
- If a folder contains 10 or more files, introduce subfolders based on semantic structure (e.g., `'format/'`, `'client/'`, `'csv/'`) to keep directories manageable.

Naming & Organization Guidelines

...

Examples

...

Output Format

You must structure your response in two clearly separated blocks, each wrapped with the appropriate tags:

`<think>`

Explain how you grouped the features into logically coherent modules with clean file and folder structure.

Describe how your choices improve clarity, minimize clutter, and reflect good design principles.

`</think>`

`<solution>`

`{ "<path to file1.py>": ["feature1", "feature2"], "<path to file2.py>": ["feature3"] } </solution>`

Prompt for Converting Subgraphs into Base Classes

You are an expert software engineer tasked with designing reusable abstractions and shared data structures for a Python codebase.

...

Base Class Output Format

You must return your design as a set of code blocks grouped by target subtree and file:

General

path/to/file.py

“python

...

“

```

...
## <Subtree Name>
### path/to/file.py
```python
...
```
</solution>

## Design Strategy
Abstractions must follow system structure and dataflow analysis, not mechanical repetition.
- Shared Data Structures: define for nodes with high out-degree (outputs consumed widely). Good candidates are feature batches, inference results, or training containers. Create a global type only when field names, data types, and usage context are stable and consistent.
- Functional Base Classes: define for nodes with high in-degree (consuming many inputs). Use when multiple modules share roles (e.g., cleaning, predicting), follow common lifecycles ('run()', 'build()', 'validate()'), or rely on similar hooks.
- Principles:
  - Avoid speculative abstractions.
  - Prefer fewer, well-justified classes (typically 1–3 globally).
  - Capture structural commonality that aids extensibility and coordination.
...
## Output Formate
Wrap your entire output in two blocks:
<think>
...
</think>
<solution>
## SubtreeA
### path/to/file.py
```python
...
```
...
</solution>

```

Prompt for Mapping Feature Paths to Interfaces

You are designing modular interfaces for a large-scale Python system. You are given repository context: overview, tree structure, skeleton, data flow, base classes, upstream interfaces, target subtree, and target file.

```

### Objective
- For each feature, define exactly one interface (function or class).
- Provide imports, signature, and detailed docstring (purpose, args, returns, assumptions).
- No implementation: use 'pass'.
- One interface per block.
### Design Guidelines
- Function: simple, atomic, stateless.
- Class: stateful, multiple methods, inherits base class, or extensible.
- Prefer fewer, well-justified abstractions.
- Group only tightly related features.
- Use global abstractions sparingly.
### Output Format
Use two blocks:
<think>
reasoning
</think>
<solution>

```



```

design_itfs_for_feature(features=["feature/path", ...]):
“python
# One interface (function or class) with docstring and pass
“
</solution>

```

B.2 Case of Built Skeleton and Designed Interfaces

We present the skeleton generated by o3-mini, together with the mapping between the generated skeleton and the nodes of machine learning algorithms. In addition, we illustrate one or two designed base classes as well as concrete functions or classes.

Repository Skeleton

```

src/
  algorithms/
    advanced/
      enhancements/
        general_enhancements/
          __init__.py
          active_learning_strategies.py
          classification_clustering_enhancements.py
          misc_general_enhancements.py
          optimization_and_meta_learning.py
          regression_enhancements.py
        __init__.py
      extended_techniques/
        extended_methods/
          __init__.py
          interpolation_and_model_learning.py
          validation_and_clustering.py
        new_models/
          __init__.py
          new_model_techniques.py
          new_model_techniques_additional.py
        __init__.py
      baselines.py
    supplemental_algorithms/
      __init__.py
      advanced_clustering_and_dimensionality_methods.py
      advanced_tokenization_and_perceptron.py
      classification_and_feature_importance_methods.py
      diverse_algorithmic_methods.py
      ensemble_evaluation_and_anomaly_detection.py
      meta_optimization_methods.py
      model_optimization_methods.py
      numerical_interpolation_methods.py
      regression_and_svm_optimization_methods.py
      spline_interpolation_and_adjusted_classifiers.py
      svm_ensemble_and_optimization_methods.py
      tokenization_methods.py
      __init__.py
    ensemble_tree/
      boosting_bagging/
        boosting/
          __init__.py
          boosting_advanced_features.py
          boosting_algorithms.py
          boosting_parameter_tuning.py
      stacking_voting/
        __init__.py
        primary.py
        secondary.py
      __init__.py
      bagging.py
      gradient_boosting.py
    decision_trees/
      __init__.py
      gradient_boosting_tree.py
      id3.py
      post_pruning.py
      random_forest.py
      regression_tree.py

```

```

__init__.py
regression/
  linear_models/
    __init__.py
    lasso.py
    multiple_linear.py
    polynomial.py
    simple_linear.py
  __init__.py
  elastic_net_regression.py
  ridge_classifier.py
supervised/
  classification/
    logistic/
      __init__.py
      cost.py
      optimization.py
      sigmoid.py
    __init__.py
    decision_tree.py
    knearest.py
    naive_bayes.py
    support_vector.py
  __init__.py
unsupervised/
  clustering/
    __init__.py
    advanced_clustering.py
    kmeans.py
    supplemental_clustering.py
  dimensionality_reduction/
    __init__.py
    extended_dr.py
    kernel_pca.py
    pca.py
  __init__.py
__init__.py
core/
  data_conversion/
    __init__.py
    api_requests.py
    feature_encoding.py
    feature_extraction.py
    format_conversion.py
    sql_queries.py
  data_transform/
    __init__.py
    filter_advanced.py
    filter_basic.py
    join_operations.py
    scaling_advanced.py
    scaling_basic.py
    sorting.py
    splitting.py
  numerics/
    __init__.py
    basic_statistics.py
    block_multiplication.py
    decompositions.py
    dot_products.py
    integration_and_distances.py
    inversions.py
    matrix_factorization.py
    matrix_rearrangements.py
    regression_statistics.py
    sparse_lu.py
  preprocessing/
    __init__.py
    csv_io.py
    data_cleaning.py
    dimensionality_analysis.py
    inverse_transformations.py
    json_io.py
    log_transformations.py
    noise_augmentation.py
  __init__.py
data_processing/
  analysis_pipeline/
    analytical/

```

```

__init__.py
aggregation_algorithms.py
data_perturbation.py
data_query.py
join_operations.py
list_manipulation.py
sample_partition.py
seasonal_analysis.py
pipeline_utilities/
__init__.py
data_streaming.py
learning_setup.py
model_validation.py
performance_metrics.py
__init__.py
cleaning_preparation/
advanced/
__init__.py
duplicate_handling.py
imputation_methods.py
outlier_detection.py
preparation/
__init__.py
data_splitting.py
imputation_labeling.py
validation.py
__init__.py
type_conversion.py
integration_merge/
__init__.py
aggregation_retrieval.py
merge_operations.py
merge_search.py
integration_storage/
__init__.py
api_operations.py
conversion_auth.py
dictionary_config.py
export_integration.py
io_logging.py
manipulation/
__init__.py
data_manipulation.py
shuffling.py
string_pivot/
__init__.py
pivoting.py
string_operations.py
transformation_feature_eng/
__init__.py
feature_extraction_encoding.py
file_io.py
text_enhancements.py
transformation_normalization.py
utilities/
__init__.py
metadata.py
metrics.py
parallel.py
sparse_storage.py
text_processing.py
validation/
inspection/
__init__.py
overview.py
sorting.py
statistics.py
__init__.py
data_integrity.py
__init__.py
extended_eval/
diagnostics/
__init__.py
model_quality.py
statistical_diagnostics.py
temporal_analysis.py
__init__.py
explainability.py
predictive_assessment.py

```

```

robustness.py
math_utils/
  algorithms/
    core_techniques/
      __init__.py
      clustering_and_detection.py
      fairness_and_feature_analysis.py
      matrix_operations.py
      optimization_and_selection.py
      statistical_methods.py
    __init__.py
  auxiliary/
    __init__.py
    data_manipulation.py
    geometric_operations.py
    gradient_and_imaging.py
    math_computations.py
    ml_utilities.py
    optimization_methods.py
    outlier_validation.py
    random_operations.py
    tensor_and_likelihood.py
    text_processing.py
    time_processing.py
  data_preprocessing/
    __init__.py
    model_persistence.py
    sampling.py
    text_tools.py
  performance/
    __init__.py
    drift_detection.py
    statistical_tests.py
    system_monitoring.py
    time_series_analysis.py
  pipeline_evaluation/
    __init__.py
    advanced_analysis.py
    data_resampling.py
    evaluation_metrics.py
    hyperparameter_tuning.py
    model_export_and_cv.py
    online_learning_support.py
    performance_benchmarking.py
    pipeline_creator.py
  simulation/
    __init__.py
    hyperparameter_tuning.py
    nearest_neighbor_search.py
    random_sampling.py
    simulation.py
    time_interpolation.py
  statistical_analysis/
    __init__.py
    descriptive_stats.py
    inferential_methods.py
    multivariate_analysis.py
    probabilistic_models.py
    survival_analysis.py
    time_series_models.py
    variance_metrics.py
  __init__.py
ml_compute/
  ml_methods/
    __init__.py
    clustering_methods.py
    dimensionality_reduction.py
    svm_validation.py
  optimization/
    __init__.py
    clustering_graph_techniques.py
    optimization_algorithms.py
    training_control.py
  stat_inference/
    __init__.py
    hypothesis_tests_advanced.py
    hypothesis_tests_basic.py
    model_evaluation_metrics.py
    model_selection_metrics.py

```

```

    parameter_estimation.py
    statistical_tests_metrics.py
time_series/
    distribution/
        __init__.py
        bayesian_likelihood.py
        distribution_estimation.py
        inferential_hypothesis.py
        inferential_variance.py
        __init__.py
    forecast_plots.py
    forecasting_algorithms.py
    __init__.py
model_eval/
    deployment/
        __init__.py
        deployment_ops.py
        testing_resources.py
    evaluation/
        diagnostics/
            __init__.py
            advanced_diagnostics.py
            basic_diagnostics.py
            performance_metrics.py
            __init__.py
        additional_analysis.py
        error_display.py
        regression_diagnostics.py
        training_ops.py
    management/
        __init__.py
        counterfactual.py
        feature_explanation.py
        hyperparameter.py
        persistence_ops.py
        pipeline_integration.py
        uncertainty.py
    visualization/
        __init__.py
        dashboard.py
        visual_reports.py
    __init__.py
pipeline/
    data_cleaning/
        __init__.py
        dimensionality.py
        encoding.py
        filtering.py
        imputation.py
        knn_methods.py
        merging.py
    deployment/
        __init__.py
        export.py
        integration_testing.py
        interactive.py
        monitoring.py
        online.py
        reporting_formats.py
        standard_plots.py
        trend_analysis.py
    evaluation/
        __init__.py
        cross_validation.py
        gd_optimization.py
        metrics_plots.py
        misc_evaluation.py
        monitoring.py
    feature_engineering/
        __init__.py
        basic_transformation.py
        extraction.py
        interaction.py
        selection.py
        synthesis.py
    orchestration/
        __init__.py
        configuration.py
        startup.py

```

```

transformers.py
workflow.py
preprocessing/
__init__.py
advanced_parsing.py
imputation.py
input_validation.py
training/
__init__.py
training_adjustments.py
training_strategies.py
tuning/
__init__.py
calibration.py
evolution.py
gaussian.py
meta.py
parzen.py
__init__.py
ui/
interactivity/
__init__.py
domain_commands.py
general_commands.py
help_support.py
navigation_actions.py
visualization/
__init__.py
dashboard.py
standard_charts.py
__init__.py
__init__.py
main.py
setup.py

```

SubGraph-to-Skeleton

```

Machine Learning Algorithms
  AdvancedExtensions [-> dir: src/algorithms/advanced]
    Miscellaneous [-> dir: src/algorithms/advanced/supplemental_algorithms]
      OtherAlgorithms [-> file: src/algorithms/advanced/supplemental_algorithms/
        classification_and_feature_importance_methods.py, file: src/algorithms/advanced/
        supplemental_algorithms/regression_and_svm_optimization_methods.py, file: src/
        algorithms/advanced/supplemental_algorithms/
        advanced_clustering_and_dimensionality_methods.py, file: src/algorithms/advanced/
        supplemental_algorithms/meta_optimization_methods.py, file: src/algorithms/
        advanced/supplemental_algorithms/svm_ensemble_and_optimization_methods.py, file:
        src/algorithms/advanced/supplemental_algorithms/
        spline_interpolation_and_adjusted_classifiers.py, file: src/algorithms/advanced/
        supplemental_algorithms/numerical_interpolation_methods.py, file: src/algorithms/
        advanced/supplemental_algorithms/diverse_algorithmic_methods.py, file: src/
        algorithms/advanced/supplemental_algorithms/advanced_tokenization_and_perceptron.
        py, file: src/algorithms/advanced/supplemental_algorithms/
        ensemble_evaluation_and_anomaly_detection.py, file: src/algorithms/advanced/
        supplemental_algorithms/tokenization_methods.py, file: src/algorithms/advanced/
        supplemental_algorithms/model_optimization_methods.py]
      ExtendedTechniques [-> dir: src/algorithms/advanced/extended_techniques]
        ExtendedMethods [-> file: src/algorithms/advanced/extended_techniques/extended_methods/
          interpolation_and_model_learning.py, file: src/algorithms/advanced/
          extended_techniques/extended_methods/validation_and_clustering.py]
        NewModels [-> file: src/algorithms/advanced/extended_techniques/new_models/
          new_model_techniques.py, file: src/algorithms/advanced/extended_techniques/
          new_models/new_model_techniques_additional.py]
        Baselines [-> file: src/algorithms/advanced/extended_techniques/baselines.py]
      EnhancementsAndGeneral [-> dir: src/algorithms/advanced/enhancements/general_enhancements]
        GeneralEnhancements [-> file: src/algorithms/advanced/enhancements/general_enhancements
          /optimization_and_meta_learning.py, file: src/algorithms/advanced/enhancements/
          general_enhancements/regression_enhancements.py, file: src/algorithms/advanced/
          enhancements/general_enhancements/active_learning_strategies.py, file: src/
          algorithms/advanced/enhancements/general_enhancements/misc_general_enhancements.py
          , file: src/algorithms/advanced/enhancements/general_enhancements/
          classification_clustering_enhancements.py]
      Regression [-> dir: src/algorithms/regression/linear_models]
        LinearModels [-> dir: src/algorithms/regression/linear_models]
          MultipleLinear [-> file: src/algorithms/regression/linear_models/multiple_linear.py,
            file: src/algorithms/regression/linear_models/lasso.py]

```

```

PolynomialRegression [-> file: src/algorithms/regression/linear_models/polynomial.py]
Lasso [-> file: src/algorithms/regression/linear_models/multiple_linear.py, file: src/
algorithms/regression/linear_models/lasso.py]
SimpleLinear [-> file: src/algorithms/regression/linear_models/simple_linear.py]
OtherRegression [-> dir: src/algorithms/regression]
RidgeRegressionClassification [-> file: src/algorithms/regression/ridge_classifier.py]
ElasticNet [-> file: src/algorithms/regression/elastic_net_regression.py]
UnsupervisedLearning [-> dir: src/algorithms/unsupervised]
DimensionalityReduction [-> dir: src/algorithms/unsupervised/dimensionality_reduction]
KernelMethods [-> file: src/algorithms/unsupervised/dimensionality_reduction/kernel_pca
.py]
OtherDR [-> file: src/algorithms/unsupervised/dimensionality_reduction/extended_dr.py]
PCA [-> file: src/algorithms/unsupervised/dimensionality_reduction/pca.py]
Clustering [-> dir: src/algorithms/unsupervised/clustering]
KMeans [-> file: src/algorithms/unsupervised/clustering/kmeans.py]
OtherClustering [-> file: src/algorithms/unsupervised/clustering/
supplemental_clustering.py]
AdvancedClustering [-> file: src/algorithms/unsupervised/clustering/advanced_clustering
.py]
EnsembleAndTreeMethods [-> dir: src/algorithms/ensemble_tree/boosting_bagging]
BoostingBagging [-> dir: src/algorithms/ensemble_tree/boosting_bagging]
StackingVoting [-> file: src/algorithms/ensemble_tree/boosting_bagging/stacking_voting/
secondary.py, file: src/algorithms/ensemble_tree/boosting_bagging/stacking_voting/
primary.py]
Bagging [-> file: src/algorithms/ensemble_tree/boosting_bagging/bagging.py]
GradientBoosting [-> file: src/algorithms/ensemble_tree/boosting_bagging/
gradient_boosting.py]
Boosting [-> file: src/algorithms/ensemble_tree/boosting_bagging/boosting/
boosting_algorithms.py, file: src/algorithms/ensemble_tree/boosting_bagging/
boosting/boosting_advanced_features.py, file: src/algorithms/ensemble_tree/
decision_trees/gradient_boosting_tree.py, file: src/algorithms/ensemble_tree/
boosting_bagging/boosting/boosting_parameter_tuning.py]
DecisionTrees [-> dir: src/algorithms/ensemble_tree/decision_trees]
ID3 [-> file: src/algorithms/ensemble_tree/decision_trees/id3.py]
RegressionTree [-> file: src/algorithms/ensemble_tree/decision_trees/regression_tree.py
]
GradientBoostingTree [-> file: src/algorithms/ensemble_tree/boosting_bagging/boosting/
boosting_advanced_features.py, file: src/algorithms/ensemble_tree/decision_trees/
gradient_boosting_tree.py]
PostPruning [-> file: src/algorithms/ensemble_tree/decision_trees/post_pruning.py]
RandomForest [-> file: src/algorithms/ensemble_tree/decision_trees/random_forest.py]
SupervisedLearning [-> dir: src/algorithms/supervised/classification/logistic]
Classification [-> dir: src/algorithms/supervised/classification/logistic]
LogisticRegression [-> file: src/algorithms/supervised/classification/logistic/sigmoid.
py, file: src/algorithms/supervised/classification/logistic/optimization.py, file:
src/algorithms/supervised/classification/logistic/cost.py]
SupportAndTree [-> file: src/algorithms/supervised/classification/decision_tree.py,
file: src/algorithms/supervised/classification/support_vector.py]
NaiveBayes [-> file: src/algorithms/supervised/classification/naive_bayes.py]
KNearest [-> file: src/algorithms/supervised/classification/knearest.py]

```

Base Classes Example

```

## General
### src/general/base_components.py
'''python
from abc import ABC, abstractmethod

class BaseComponent(ABC):
    """
    Abstract base class for core pipeline components.

    Provides a standard lifecycle:
    - initialize(): Prepare the component.
    - process(data): Process input data.
    - finalize(): Clean up resources.

    Concrete subclasses must implement process().
    """

    def __init__(self) -> None:
        self.initialized: bool = False

    def initialize(self) -> None:
        """
        Set up the component before processing starts.

```



```

        """
        self.initialized = True

    @abstractmethod
    def process(self, data):
        """
        Process the input data.

        Args:
            data: Input data in a predefined schema.

        Returns:
            Output data after processing.
        """
        pass

    def finalize(self) -> None:
        """
        Tear down or clean up the component after processing.
        """
        self.initialized = False

class EstimatorComponent(BaseComponent):
    """
    Abstract base class for estimator components (e.g., models).

    Defines the contract for model training and prediction.
    """

    @abstractmethod
    def fit(self, training_data) -> None:
        """
        Train the estimator using the provided training data.

        Args:
            training_data: Data batch conforming to a shared TrainingBatch schema.
        """
        pass

    @abstractmethod
    def predict(self, input_data):
        """
        Generate predictions based on the input data.

        Args:
            input_data: Data in a format specified by the pipeline.

        Returns:
            Predictions corresponding to the input features.
        """
        pass

    def process(self, data):
        """
        For an estimator, process() defaults to prediction.
        """
        return self.predict(data)
'''

```

Designed Interfaces

```

## Feature Paths: "Regression/LinearModels/PolynomialRegression/Cubic Regression with
regularization", "Regression/LinearModels/PolynomialRegression/polynomial model fitting", "
Regression/LinearModels/PolynomialRegression/cubic regression", "Regression/LinearModels/
PolynomialRegression/quadratic regression", "Regression/LinearModels/PolynomialRegression/
Quadratic Regression with regularization"
### src/algorithms/regression/linear_models/polynomial.py
from src.general.base_components import EstimatorComponent
class PolynomialRegressor(EstimatorComponent):

    def __init__(self, degree: int, regularization_lambda: float=0.0) -> None:
        pass

    def fit(self, X: list[float], y: list[float]) -> None:
        """
        Fit the polynomial regression model to the provided data.

```

```

Constructs the polynomial features based on the specified degree and applies
optional regularization if regularization_lambda is provided (> 0).

Args:
    X (list[float]): A list of feature values.
    y (list[float]): A list of target values corresponding to the features.

Returns:
    None

Raises:
    ValueError: If the degree is not supported or if input lists are empty or mismatched.
    """
    pass

def predict(self, X: list[float]) -> list[float]:
    """
    Generate predictions using the fitted polynomial regression model.

    Transforms the input features into polynomial features and computes
    the output via the fitted model coefficients. Applies regularization adjustments
    if the model was fitted with a regularization term.

    Args:
        X (list[float]): A list of feature values for prediction.

    Returns:
        list[float]: A list of predicted values.
    """
    pass

```

B.3 Patterns in Implementation-Level Graph Construction

The mapping from RPGs to code structures exhibits a strong isomorphic relationship: each subgraph corresponds to a coherent code region, with files, classes, and functions serving as structural anchors. Table 7 illustrates this correspondence for the case of *o3-mini* during *sklearn* generation, where algorithmic subgraphs (e.g., ML Algorithms, Data Processing, ML Pipeline) map to a larger number of files and functions, while auxiliary subgraphs (e.g., Diagnostics, Visualization) remain compact yet feature-dense. This pattern reflects the semantic granularity of different subgraphs: core computational domains require broader structural scaffolding, whereas specialized domains concentrate more features per unit. Extending to the cross-repository view in Table 6, we observe that both models preserve this structural isomorphism but with distinct emphases: *o3-mini* tends to distribute features more evenly across units, while *qwen3-coder* consistently produces the highest feature densities, especially at the class level. Together, these results demonstrate that the graph-to-code translation process not only preserves the hierarchical semantics of the RPG but also manifests in distinct structural footprints that vary with model choice.

Table 6: Per-repository structural statistics across *o3-mini* and *qwen-coder*. “Count” = number of entities (Files/Classes/Functions) per repository; “Avg Feat.” = mean number of features per entity (features per file/class/function).

| Repo | o3-mini | | | | | | qwen-coder | | | | | |
|----------------|---------|-----------|---------|-----------|-----------|-----------|------------|-----------|---------|-----------|-----------|-----------|
| | Files | | Classes | | Functions | | Files | | Classes | | Functions | |
| | Count | Avg Feat. | Count | Avg Feat. | Count | Avg Feat. | Count | Avg Feat. | Count | Avg Feat. | Count | Avg Feat. |
| TableKit | 475 | 3.64 | 252 | 2.28 | 1092 | 1.05 | 271 | 6.69 | 496 | 1.62 | 587 | 1.50 |
| MLKit-Py | 266 | 4.74 | 321 | 1.64 | 708 | 1.04 | 566 | 2.44 | 815 | 1.30 | 281 | 1.13 |
| StatModeler | 219 | 4.71 | 117 | 2.47 | 726 | 1.02 | 330 | 3.48 | 573 | 1.22 | 411 | 1.07 |
| SymbolicMath | 126 | 4.70 | 95 | 2.17 | 370 | 1.04 | 89 | 8.98 | 71 | 1.73 | 786 | 0.86 |
| PyWebEngine | 440 | 3.89 | 576 | 1.74 | 689 | 1.02 | 482 | 3.52 | 890 | 1.26 | 501 | 1.13 |
| HttpEasy | 104 | 4.17 | 79 | 2.43 | 235 | 1.03 | 178 | 4.28 | 239 | 1.52 | 366 | 1.06 |
| Average | 271.7 | 4.31 | 240 | 2.12 | 636.7 | 1.03 | 319.3 | 4.90 | 514 | 1.44 | 488.7 | 1.12 |

Table 7: Structural distribution of files, classes, functions, and feature densities corresponding to each **subgraph in the feature graph of o3-mini during sklearn synthesis**. Here, “Files/Classes/Functions” denote the number of code units mapped from each subgraph; “File/Class/Function Features” are the total extracted features; and “Avg Features/...” indicates the average number of features per unit type.

| Subgraph | Files | Classes | Functions | File Features | Class Features | Function Features | Avg Feat./File | Avg Feat./Class | Avg Feat./Func |
|-----------------|-------|---------|-----------|---------------|----------------|-------------------|----------------|-----------------|----------------|
| ML Algorithms | 58 | 171 | 67 | 323 | 256 | 67 | 5.57 | 1.50 | 1.00 |
| Math Utilities | 47 | 26 | 102 | 143 | 40 | 103 | 3.04 | 1.54 | 1.01 |
| Data Processing | 45 | 30 | 169 | 231 | 59 | 172 | 5.13 | 1.97 | 1.02 |
| ML Pipeline | 38 | 23 | 149 | 202 | 42 | 160 | 5.32 | 1.83 | 1.07 |
| Core Operations | 30 | 15 | 76 | 124 | 33 | 91 | 4.13 | 2.20 | 1.20 |
| ML Computation | 19 | 39 | 34 | 88 | 54 | 34 | 4.63 | 1.38 | 1.00 |
| ML Evaluation | 17 | 12 | 51 | 77 | 25 | 52 | 4.53 | 2.08 | 1.02 |
| ML Diagnostics | 6 | 3 | 44 | 50 | 6 | 44 | 8.33 | 2.00 | 1.00 |
| Visualization | 6 | 2 | 16 | 24 | 8 | 16 | 4.00 | 4.00 | 1.00 |

C Appendix of Graph-Guided Repository Generation

C.1 Details on Localization

To facilitate the localization stage in graph-guided repository generation, we designed a graph-guided toolset that allows agents to systematically explore and map design-level features onto concrete code artifacts. The tools support both fine-grained inspection of files and interfaces, as well as feature-driven exploration across the repository. Specifically, `view_file_interface_feature_map` and `get_interface_content` enable inspection of code structures and retrieval of their implementations, while `expand_leaf_node_info` and `search_interface_by_functionality` allow navigation of the RPG and fuzzy semantic search. Finally, the `Terminate` command ensures that the localization process produces a ranked and standardized output. Together, these tools provide a structured workflow that balances automation with flexibility, ensuring both accuracy and interpretability in the localization process.

Localization Tools

```

### Interface Inspection Tools
- 'view_file_interface_feature_map(file_path)'
  Inspects a single Python file to list the interface structures (functions, classes, methods) it
  contains, along with the feature mappings they support.
  *Usage*: Useful for quickly understanding which interfaces exist in a given file and the feature
  tags associated with them.
  *Example*:
  \begin{verbatim}
  view_file_interface_feature_map('src/algorithms/classifier.py')
  \end{verbatim}

- 'get_interface_content(target_specs)'
  Retrieves the full implementation code of a specific function, class, or method, given its fully
  qualified name (file path + entity name).
  *Usage*: Applied when a particular interface has been located and its source code needs to be
  examined in detail.
  *Example*:
  \begin{verbatim}
  get_interface_content(['src/core/data_loader.py:DataLoader.load_data'])
  get_interface_content(['src/core/utils.py:clean_text'])
  \end{verbatim}

### Feature-Driven Exploration Tools
- 'expand_leaf_node_info(feature_path)'
  Given a feature path from the implemented feature tree, this tool expands and lists all
  associated interfaces (functions or classes) in a structural summary.
  *Usage*: Applied when analyzing how a specific functional leaf node in the design tree maps to
  repository interfaces.
  *Example*:
  \begin{verbatim}
  expand_leaf_node_info('Algorithm/Supervised Learning/Classification Algorithms/Naive Bayes')
  \end{verbatim}

- 'search_interface_by_functionality(keywords)'
  Performs a fuzzy semantic search for interfaces based on given keywords and returns the top-5
  most relevant interface implementations.

```

```

*Usage*: Useful when the exact file or interface name is unknown, but functionality-related
keywords are available.
*Example*:
\begin{verbatim}
search_interface_by_functionality(['optimize', 'initialize'])
\end{verbatim}

### Termination Tool
- 'Terminate(result)'
  Terminates the localization exploration and returns the final ranked list of located interfaces.
  The result must follow the specified JSON-style format, including file path and interface
  type (function, class, or method).
*Usage*: Invoked after completing exploration to deliver the final interface localization results
.
*Example*:
\begin{verbatim}
Terminate(result=[
  {"file_path": "top1_file_fullpath.py", "interface": "method: Class1.function1"},
  {"file_path": "top2_file_fullpath.py", "interface": "function: function2"},
  {"file_path": "top3_file_fullpath.py", "interface": "class: Class3"},
])
\end{verbatim}

```

C.2 Tools for Coding

To support systematic repository modification, we introduce a suite of editing tools that operate at multiple levels of code granularity.

Editing Coding Tools

```

### edit_whole_class_in_file(file_path, class_name)
Use when: An entire class (including all methods and its docstring) requires editing or replacement
.
Output must: Provide the full class definition, with all methods and docstring.

### edit_method_of_class_in_file(file_path, class_name, method_name)
Use when: A single method inside a class needs to be edited or replaced.
Output must:
- Return the full 'class ClassName:' block containing only the target method.
- Exclude all unrelated methods.
- Do not output the method alone; it must appear within its class block.

### edit_function_in_file(file_path, function_name)
Use when: A top-level function requires editing or replacement.
Output must: Provide the full function, including signature, logic, and docstring.

### edit_imports_and_assignments_in_file(file_path)
Use when: Import statements or top-level assignments in the file need to be added or corrected.
Output must:
- Contain only import statements and top-level assignments (no functions or classes).
- Follow a strict import order: (1) standard library, (2) third-party packages, (3) local modules.
Important:
- Do not remove existing imports unless they are demonstrably incorrect (e.g., invalid path, typo,
  or non-existent module).
- Retain imports even if they appear unused, to preserve runtime dependencies.

### Terminate()
Use when: All required edits are complete and no further tool invocations are necessary.
Important: Do not call prematurely; ensure the issue is fully resolved before termination.

```

C.3 Graph-Guided Localization and Editing Logs

To illustrate how the proposed system operates in practice, we include representative logs that demonstrate the end-to-end process: (i) locating relevant interfaces in the repository through RPG-guided exploration, and (ii) applying code editing operations to implement the required functionality. These logs capture the step-by-step decision process, tool invocations, and final termination results.

Task Example Trajectory

Task: This issue is to implement the `AdvancedDifferentialTransformer` class in `src/symbolic/enhanced/calculus_summation/differential.py`. Contributors should aim to encapsulate the necessary transformation logic and ensure the methods integrate well with the rest of the symbolic algebra framework. Clear error handling and input validation are expected in order to gracefully manage undefined behaviors.

=====Localization=====

Agent: <think>

I will begin by exploring the file "`src/symbolic/enhanced/calculus_summation/differential.py`" since the `AdvancedDifferentialTransformer` class should be implemented here. I will use the `view_file_interface_feature_map` tool on that file to inspect its structure and get an idea of what is already provided, and any hints or TODOs might be present.

</think>

<solution>

```
'''
view_file_interface_feature_map('src/symbolic/enhanced/calculus_summation/differential.py')
'''
```

</solution>

Env: '''python

```
## src/symbolic/enhanced/calculus_summation/differential.py
```

```
from typing import Any, Dict
```

```
### Features: Calculus & Advanced Summation/Differential & Asymptotic Methods/Advanced
Transformation Methods/asymptotic series expansion, Calculus & Advanced Summation/Differential
& Asymptotic Methods/Advanced Transformation Methods/differentiate special functions,
Calculus & Advanced Summation/Differential & Asymptotic Methods/Advanced Transformation
Methods/simplify inequality expressions
class AdvancedDifferentialTransformer:
```

```
    def asymptotic_series_expansion(self, expression: Any, expansion_order: int) -> Any:
        ...
```

```
    def differentiate_special_function(self, expression: Any, function_name: str, additional_params
        : Dict[str, Any]=None) -> Any:
        ...
```

```
    def simplify_inequality(self, inequality_expr: Any, tolerance: float=1e-05) -> Any:
```

```
'''
```

Agent: <think>

I'll now check the file "`src/symbolic/enhanced/simplification/algebraic.py`" since it is part of the Enhanced Symbolic Algorithms simplification area, and it might contain functions or methods relevant to simplifying inequalities.

</think>

<solution>

```
'''
view_file_interface_feature_map('src/symbolic/enhanced/simplification/algebraic.py')
'''
```

</solution>

Env: '''python

```
## src/symbolic/enhanced/simplification/algebraic.py
```

```
from typing import Any
```

```
from src.expr_transform.simplification.techniques import apply_commutative_associative_rewriting,
    order_terms
```

```
from src.symbolic.core.expression_manipulation.tree_utils import collect_like_terms
```

```
from src.symbolic.core.expression_manipulation.rewriting import canonicalize_expression_tree
```

```
### Features: Simplification/Algebraic Simplification/Basic Techniques/simplify expressions
```

```
def simplify_expression(expr: Any) -> Any:
```

```
'''
    ...
'''
```

Agent: <think>

I will now search for interfaces related to "differentiate special function" or advanced symbolic differentiation since our method `differentiate_special_function` might rely on those routines. This might reveal supporting functions or logic that can be reused.

</think>

<solution>

```
'''
search_interface_by_functionality(['differentiate special function', 'advanced symbolic
differentiation'])
'''
```

</solution>

```

Env: Keyword "differentiate special function" matched interfaces: class AdvancedExpressionAnalyzer
, class AutomaticDifferentiator, class MultivariableDifferentiator.
Keyword "advanced symbolic differentiation" matched interfaces: class AutomaticDifferentiator,
class IntegerAdditionOperations, class MultivariableDifferentiator.

## File: src/calculus/differentiation/partial_derivatives.py

from typing import Any, List, Callable

### Features: Differentiation/Specialized/Techniques/validate partial derivatives, Differentiation/
Specialized/Techniques/compute partial derivatives, Differentiation/Specialized/Techniques/
symbolic gradient computation, Differentiation/Specialized/Techniques/Jacobian matrix,
Differentiation/Specialized/Techniques/symbolic hessian computation, Differentiation/
Specialized/Techniques/multivariable differentiation
class MultivariableDifferentiator:

    def validate_partial_derivatives(self, function: Callable[..., Any], variables: List[str],
        computed: Any) -> bool:
        ...

    def compute_partial_derivative(self, function: Any, variable: str) -> Any:
        """
        ...

    def compute_symbolic_gradient(self, function: Any, variables: List[str]) -> Any:
        ...

    def compute_jacobian_matrix(self, functions: List[Any], variables: List[str]) -> Any:
        ...

    def compute_symbolic_hessian(self, function: Any, variables: List[str]) -> Any:
        ...

    def differentiate_multivariable(self, function: Any, variables: List[str]) -> Any:
        ...
...

Agent: <think>
Next, I'll search for interfaces that mention "asymptotic series expansion" because that
functionality is required in AdvancedDifferentialTransformer.asymptotic_series_expansion. This
might reveal candidate functions that can be reused to generate series expansions.

I'll also search for interfaces related to "simplify inequality" to see if there's any existing
logic for inequality simplification that we can leverage.

Let's perform both searches.
</think>
<solution>
'''
search_interface_by_functionality(['asymptotic series expansion'])
search_interface_by_functionality(['simplify inequality'])
'''
</solution>

Env: Keyword "asymptotic series expansion" matched interfaces: class AdvancedExpressionAnalyzer,
class MultivariableDifferentiator, class NumericalIntegrator.

## File: src/calculus/differentiation/partial_derivatives.py

from typing import Any, List, Callable

### Features: Differentiation/Specialized/Techniques/validate partial derivatives, Differentiation/
Specialized/Techniques/compute partial derivatives, Differentiation/Specialized/Techniques/
symbolic gradient computation, Differentiation/Specialized/Techniques/Jacobian matrix,
Differentiation/Specialized/Techniques/symbolic hessian computation, Differentiation/
Specialized/Techniques/multivariable differentiation
class MultivariableDifferentiator:
    ...

Agent: <think>
We have gathered several candidate supporting interfaces needed for our implementation of
AdvancedDifferentialTransformer. Our analysis indicates that the following components are very
likely to be relevant dependencies:
...

These are our top-5 interface-level code units that are most likely to serve as essential
dependencies for implementing the class AdvancedDifferentialTransformer with its three key
methods.

```

```

</think>
<solution>
'''
Terminate(result=[
    {"file_path": "src/symbolic/core/calculus/differentiation.py", "interface": "function:
        differentiate_expression_tree"},
    {"file_path": "src/symbolic/enhanced/simplification/algebraic.py", "interface": "function:
        simplify_expression"},
    {"file_path": "src/expr_transform/transformation/analysis_expression.py", "interface": "class:
        AdvancedExpressionAnalyzer"},
    {"file_path": "src/arith_algebra/differentiation.py", "interface": "class:
        AutomaticDifferentiator"},
    {"file_path": "src/arith_algebra/equation_solving/nonlinear_symbolic.py", "interface": "class:
        SymbolicEquationSolver"}
])
'''

=====Coding=====
Agent: <think>
We need to implement the AdvancedDifferentialTransformer class in the file src/symbolic/enhanced/
calculus_summation/differential.py. The class should provide three methods. For
asymptotic_series_expansion, we will check that expansion_order is a positive integer and that
the expression is provided. Our implementation will simulate an asymptotic expansion by
returning a new symbolic expression with a type marker (e.g., "asymptotic_series") including
the original expression and the specified order. Similarly, differentiate_special_function
should validate that the provided function name is among a set of allowed special functions (
for example, erf, gamma, beta, bessell, hypergeometric) and then simulate a specialized
differentiation by returning a new symbolic expression that indicates that the transformation
has been applied. Finally, simplify_inequality will check that the input is structured (i.e. a
dictionary) and that tolerance is a positive float reasonably bounded (we can enforce
tolerance > 0 and maybe below 1 to avoid very high tolerance) and then simulate a
simplification by deep-copying the input expression with additional metadata fields. We also
include proper error handling and input validation.
</think>
<solution>
edit_whole_class_in_file("src/symbolic/enhanced/calculus_summation/differential.py", "
    AdvancedDifferentialTransformer")
'''python
from typing import Any, Dict
import copy

class AdvancedDifferentialTransformer:
    """
    Provides advanced differential transformation methods for asymptotic series expansion,
    differentiating special functions, and simplifying inequality expressions within symbolic
    mathematics.

    This class encapsulates transformation techniques that are applied to symbolic expressions
    in the context of differential and asymptotic methods, enabling enhanced symbolic analyses and
    computations.

    Methods:
        asymptotic_series_expansion(expression: Any, expansion_order: int) -> Any:
            Expand the provided expression into its asymptotic series representation up to the
            given order.

        differentiate_special_function(expression: Any, function_name: str, additional_params: Dict
            [str, Any] = None) -> Any:
            Differentiate a symbolic expression that contains a special function using advanced
            rules.

        simplify_inequality(inequality_expr: Any, tolerance: float = 1e-5) -> Any:
            Simplify a symbolic inequality expression using transformation techniques to reduce
            complexity.
    """
    ...
'''
</solution>

Env:
# Code Patch Diff:
--- src/symbolic/enhanced/calculus_summation/differential.py
+++ src/symbolic/enhanced/calculus_summation/differential.py
@@ -1,4 +1,5 @@
from typing import Any, Dict
+import copy

class AdvancedDifferentialTransformer:
    """

```



```

@@ -18,6 +19,7 @@
    simplify_inequality(inequality_expr: Any, tolerance: float = 1e-5) -> Any:
        Simplify a symbolic inequality expression using transformation techniques to reduce
        complexity.
    """
+    ALLOWED_SPECIAL_FUNCTIONS = {'erf', 'gamma', 'beta', 'bessel', 'hypergeometric'}

    def asymptotic_series_expansion(self, expression: Any, expansion_order: int) -> Any:
        """
@@ -26,15 +28,23 @@
    Args:
        expression (Any): The symbolic expression to be expanded.
        expansion_order (int): The order up to which the series expansion should be computed.
-
+
    Returns:
        Any: A symbolic expression representing the asymptotic series expansion.

-    Edge Cases:
-        - If expansion_order is non-positive, the behavior is undefined.
-        - The function expects a valid symbolic expression input.
+    Raises:
+        ValueError: If expansion_order is not a positive integer.
+        TypeError: If the expression is not in a valid symbolic format (e.g., a dictionary).
    """
-    pass
+    if not isinstance(expansion_order, int) or expansion_order <= 0:
+    ...

Agent: <solution>
Terminate()
</solution>

```

C.4 Testing Framework

As illustrated in Algorithm 3, the testing framework combines **unit tests**, **regression tests**, and **integration tests** into a unified process. Unit tests validate each new patch at a fine-grained level, regression tests re-execute existing nodes to ensure backward compatibility, and integration tests verify that multiple patched components interact correctly when combined.

As illustrated in Algorithm 4, the testing pipeline proceeds in a sequence of stages: branch planning, test generation, execution, judgment, and repair. First, a candidate test branch is created for the given code unit(s). Then, test code is generated and wrapped into a `TestNode` or `IntegrationTestNode`, which is executed inside a controlled Docker environment. The execution results are judged by an LLM; if failures are detected, the framework automatically generates fix queries and iteratively repairs the test until a validated version is obtained.

C.5 Statistics of Three Stage

Table 8 demonstrates that graph-guided localization provides reasonable efficiency across repositories, with *Incremental Development* generally easier to localize than *Integration Testing* or *Debugging*. In terms of models, o3-mini achieves higher localization efficiency but with larger variance, whereas qwen3-coder shows more stable yet overall lower efficiency. These results suggest that while graph guidance is effective, model capacity and stability jointly influence localization performance.

Table 8: Localization results across six open-source repositories under three task categories: *Integration Testing*, *Incremental Development*, and *Debugging*. Each entry reports the mean performance with standard deviation (mean \pm std) of the corresponding model-task pair.

| Model | Task | TableKit | MLKit-Py | HttpEasy | PyWebEngine | StatModeler | SymbolicMath |
|-------------|-------------------------|------------------|-----------------|------------------|------------------|------------------|-----------------|
| o3-mini | Integration Testing | 13.33 \pm 2.92 | 8.75 \pm 4.32 | 10.94 \pm 3.44 | 6.65 \pm 1.98 | 9.24 \pm 3.65 | 7.88 \pm 3.30 |
| | Incremental Development | 12.30 \pm 5.19 | 9.83 \pm 4.00 | 11.60 \pm 5.09 | 12.51 \pm 6.67 | 12.62 \pm 6.15 | 9.93 \pm 5.13 |
| | Debugging | 11.59 \pm 5.74 | 8.24 \pm 4.40 | 9.15 \pm 5.55 | 10.28 \pm 8.50 | 13.02 \pm 7.01 | 8.90 \pm 6.21 |
| qwen3-coder | Integration Testing | 6.16 \pm 2.37 | 6.62 \pm 2.12 | 7.89 \pm 2.42 | 5.93 \pm 2.06 | 9.24 \pm 3.65 | 7.88 \pm 3.30 |
| | Incremental Development | 6.81 \pm 1.87 | 7.10 \pm 1.98 | 7.48 \pm 1.85 | 6.98 \pm 1.92 | 6.49 \pm 1.79 | 7.12 \pm 1.77 |
| | Debugging | 6.75 \pm 2.21 | 6.01 \pm 2.16 | 6.25 \pm 1.82 | 6.62 \pm 2.47 | 5.94 \pm 2.19 | 6.42 \pm 1.94 |

Algorithm 3 Patch-Oriented Testing with Unit, Regression, and Integration Stages

Require: Patch set \mathcal{P} ; repo skeleton \mathcal{R} ; dependency code D ; existing unit nodes \mathcal{N}_u ; existing integration nodes \mathcal{N}_i ; task description Θ

```

1: function TESTPATCHES( $\mathcal{P}, \mathcal{R}, D, \Theta$ )
2:    $\mathcal{T}_{unit} \leftarrow []$ ;  $\mathcal{T}_{inte} \leftarrow []$ 
3:    $\mathcal{T}_{traj} \leftarrow \{\text{unit} : \{\}, \text{inte} : \{\}\}$ 
4:    $\mathcal{P}' \leftarrow \mathcal{P} \cup \text{FINDDEPPATCHES}(\mathcal{P})$  ▷ Extend patch set with dependency patches
5:   for patch  $p \in \mathcal{P}'$  do
6:      $n_{old} \leftarrow \text{FINDEXISTINGUNITNODE}(\mathcal{N}_u, p)$ 
7:     if  $n_{old} \neq \emptyset$  and  $\text{SAMESIGNATUREORLOGIC}(n_{old}, p)$  then
8:        $n_{new} \leftarrow n_{old}$  ▷ Regression test: reuse existing node if signature/logic unchanged
9:     else
10:       $n_{new}, traj \leftarrow \text{CREATEORUPDATEUNITNODE}(p, D, \Theta, n_{old})$ 
11:       $\mathcal{T}_{traj}[\text{unit}][p.\text{key}] \leftarrow traj$ 
12:    end if
13:     $\mathcal{R}.\text{INSERTFILE}(n_{new}.\text{test\_file}, n_{new}.\text{test\_code})$ 
14:     $res \leftarrow n_{new}.\text{EXECUTETEST}()$ 
15:     $\mathcal{T}_{unit}.\text{APPEND}(res)$ 
16:  end for
17:  for patch group  $\mathcal{G}$  clustered by integration-node do
18:     $n_{old} \leftarrow \text{FINDEXISTINGINTEGRATIONNODE}(\mathcal{N}_i, \mathcal{G})$ 
19:    if  $n_{old} \neq \emptyset$  and  $\text{ALLEQUAL}(n_{old}, \mathcal{G})$  then
20:       $n_{new} \leftarrow n_{old}$  ▷ Regression integration test: reuse existing node
21:    else
22:       $n_{new}, traj \leftarrow \text{CREATEINTEGRATIONNODE}(\mathcal{G}, \Theta)$ 
23:       $\mathcal{T}_{traj}[\text{inte}][\mathcal{G}] \leftarrow traj$ 
24:       $\mathcal{R}.\text{INSERTFILE}(n_{new}.\text{test\_file}, n_{new}.\text{test\_code})$ 
25:    end if
26:     $res \leftarrow n_{new}.\text{EXECUTETEST}()$ 
27:     $\mathcal{T}_{inte}.\text{APPEND}(res)$ 
28:  end for
29:  return  $\mathcal{T}_{unit} \cup \mathcal{T}_{inte}, \mathcal{T}_{traj}$ 
30: end function

```

As shown in Table 9, **o3-mini** achieves relatively high code success rates across repositories, often exceeding 75% and in some cases approaching 90%, whereas **qwen3-coder** lags behind with rates around 50–55%. In contrast, the corresponding test coverage remains moderate, typically within the 60–70% range. Figure 12 further illustrates that coverage fluctuates and tends to decline as code length increases: shorter implementations reach high class-level coverage, but both function-level and overall coverage drop significantly with greater complexity. These results suggest that while current models are increasingly effective at generating functional code, their ability to produce comprehensive and high-quality test cases remains limited, highlighting test generation as a key bottleneck for practical deployment.

Table 9: Average success rate and test coverage (%) for six repositories across two models.

| Model | TableKit | | MLKit-Py | | HttpEasy | | PyWebEngine | | StatModeler | | SymbolicMath | |
|-------------|----------|----------|----------|----------|----------|----------|-------------|----------|-------------|----------|--------------|----------|
| | Success | Coverage | Success | Coverage | Success | Coverage | Success | Coverage | Success | Coverage | Success | Coverage |
| o3-mini | 81.8% | 65.0% | 82.8% | 61.0% | 88.9% | 64.0% | 74.7% | 60.0% | 71.0% | 62.0% | 84.8% | 59.0% |
| qwen3-coder | 55.0% | 48.0% | 52.0% | 46.0% | 50.0% | 45.0% | 53.0% | 47.0% | 54.0% | 48.0% | 51.0% | 46.0% |

D Details about RepoCraft Benchmark

In this section, we describe the construction of the REPOCRAFT benchmark, covering four key aspects: the choice of repositories, the preparation of test data, the evaluation methodology, and the configuration of agent systems.

Algorithm 4 End-to-End Test Generation, Execution, and Repair

Require: Repo skeleton \mathcal{R} ; tested unit(s) U ; source code C ; optional prior test node n_{old} ; maximum retries T_{\max}

```

1: function RUNTESTINGPIPELINE( $\mathcal{R}, U, C, n_{old}$ ) ▷ Main entry point for testing workflow
2:   // — Step 1: Plan test branches —
3:    $branch \leftarrow \text{GENERATECODEBRANCH}(C, n_{old}, T_{\max})$ 
4:   // — Step 2: Generate candidate test code —
5:    $test\_code \leftarrow \text{GENERATETEST}(branch, C, U, n_{old})$ 
6:   // — Step 3: Build a TestNode —
7:   if  $U$  represents integration of multiple units then
8:      $n \leftarrow \text{INTEGRATIONTESTNODE}(U, test\_code)$ 
9:   else
10:     $n \leftarrow \text{UNITTESTNODE}(U, C, test\_code)$ 
11:   end if
12:   // — Step 4: Execute test code in Docker —
13:    $result \leftarrow n.\text{EXECUTETEST}()$ 
14:    $output \leftarrow result.stdout \parallel result.stderr$ 
15:   // — Step 5: LLM judge outcome —
16:   if  $result$  contains errors then
17:      $(err\_type, reviews) \leftarrow \text{LLMJUDGE}(C, test\_code, output, branch)$ 
18:     if  $err\_type \in \{\text{test\_code}, \text{environment}\}$  then
19:        $query \leftarrow \text{GENERATEFIXQUERY}(C, test\_code, output, branch, reviews)$ 
20:        $n \leftarrow \text{FIXTESTANDENV}(query, U, C, output, n)$ 
21:     end if
22:   end if
23:   // — Step 7: Return final validated test node —
24:   return  $n$ 
25: end function

```

D.1 Repositories Selection

For the benchmark, we curated six representative open-source repositories: *scikit-learn*, *pandas*, *Django*, *statsmodels*, *SymPy*, and *requests*. These projects span diverse functional domains including machine learning, data analysis, web frameworks, statistical modeling, symbolic computation, and HTTP communication, thereby ensuring broad coverage of typical software development tasks. To prevent models from simply memorizing or retrieving solutions from training data, we deliberately anonymized the repositories by modifying their names and descriptions. Furthermore, the task instructions prohibit directly reusing the original implementations, requiring models to generate solutions guided only by feature specifications. This setup enforces a fairer evaluation, focusing on the models’ capacity for feature-grounded reasoning and code generation rather than exploitation of prior exposure.

D.2 Evaluation Tasks Collection

To construct a diverse and reliable evaluation set, we developed an automated pipeline that extends and systematizes the collection of test functions from the official repositories. Our design leverages the fact that mature open-source projects typically include comprehensive test suites with robust inputs and ground-truth outputs, ranging from unit-level checks to integration-level workflows. These tests provide a principled source of evaluation data, ensuring that generated repositories are assessed on both algorithmic diversity and functional correctness.

Test Function Harvesting. For each repository, we first gathered all available test functions and classes. These serve as the raw pool of evaluation candidates, capturing the behaviors developers themselves deemed important to verify.

Hierarchical Categorization. Next, we organized the collected tests into a hierarchical taxonomy. At the top level, categories follow the natural modular structure used by human developers (e.g., *metrics*, *linear_model*, *decomposition*). Within each category, we grouped related test classes and functions by algorithmic target. For example:

```

{
  "metrics": {
    "test_regression": {
      "functions": {

```

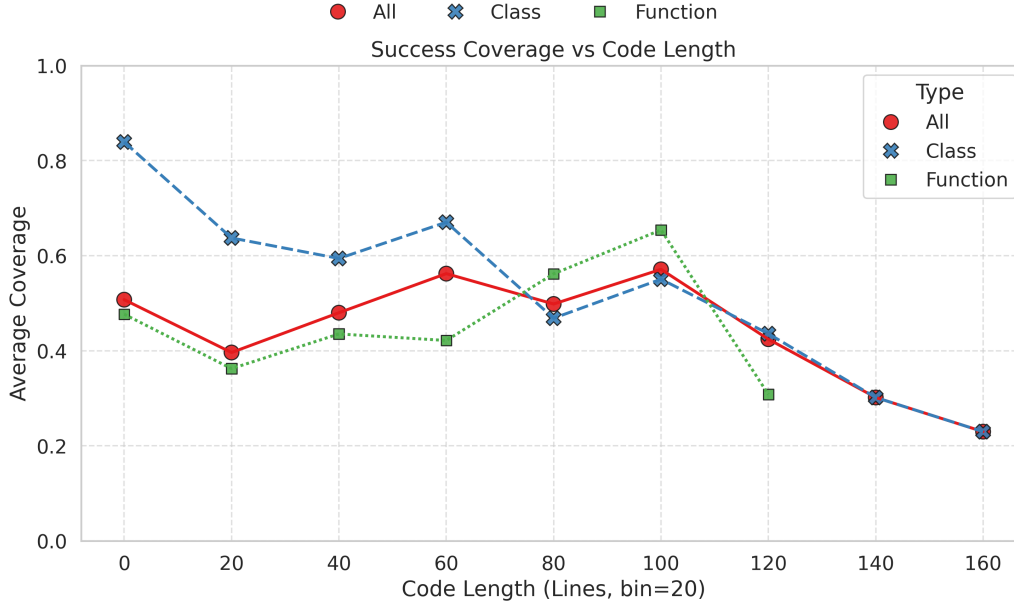


Figure 12: Test coverage of **o3-mini** on **MLKit-Py** during generation. The figure shows how the coverage of generated test functions varies as code length increases.

```

"reg_targets": [
  "test__check_reg_targets",
  "test__check_reg_targets_exception"
],
"regression_metrics": [
  "test_regression_metrics",
  "test_root_mean_squared_error_multioutput_raw_value",
  ...
],
"pinball_loss": [
  "test_mean_pinball_loss_on_constant_predictions",
  "test_dummy_quantile_parameter_tuning",
  "test_pinball_loss_relation_with_mae"
]
}
}
}
}

```

This taxonomy mirrors repository semantics: higher levels correspond to broad functional modules, while deeper levels capture fine-grained algorithmic tests.

Sampling and Filtering. To ensure balanced coverage, we applied the sampling algorithm(Alg 1) to draw representative subsets of test categories. Each sampled test was then refined into a task description that models could follow during generation. Finally, we filtered out cases irrelevant to core algorithmic behavior (e.g., string formatting checks, version consistency tests), retaining only tests that probe substantive computational functionality.

Example Task Instance. To illustrate the outcome of the pipeline, consider the following task specification extracted from the *Django* repository:

```

{
  "category": "gis_migrations",
  "file": "tests/gis_tests/gis_migrations/test_operations.py",
  "module": "class OperationTests",
  "cap": "spatial_index",
  "functions": [
    "test_create_model_spatial_index",
    "test_alter_field_add_spatial_index",
    "test_alter_field_remove_spatial_index",
    "test_alter_field_nullable_with_spatial_index",
    "test_alter_field_with_spatial_index"
  ]
}

```

Table 10: Overview of the six benchmark repositories in REPOCRAFT. Each repository is anonymized by renaming to prevent direct memorization or retrieval by models. We list both the anonymized names and their original counterparts, together with category, purpose, and scope.

| Original Name | Anonymized Name | Category | Purpose | Scope |
|---------------|-----------------|------------------------------|--|--|
| scikit-learn | MLKit-Py | Machine Learning Framework | Provides efficient tools for data mining and analysis, supporting classical supervised and unsupervised learning algorithms. | Focuses on model training, evaluation, and selection for standard ML tasks; excludes deep learning and distributed systems. |
| sympy | SymbolicMath | Symbolic Computation Library | Enables symbolic mathematics including algebraic manipulation, calculus, simplification, and equation solving. | Pure Python implementation of symbolic and algebraic computation, lightweight and extensible, without external dependencies. |
| pandas | TableKit | Data Analysis Library | Provides flexible data structures (e.g., DataFrame, Series) for manipulating and analyzing tabular data. | Supports efficient single-machine structured data processing; excludes distributed data frameworks. |
| django | PyWebEngine | Web Framework | High-level framework for rapid development with ORM, routing, templating, and admin support. | Offers an all-in-one toolkit for building web applications on small-scale/single-server deployments. |
| requests | HttpEasy | HTTP Client Library | Simple, human-friendly library for sending HTTP requests and handling responses. | Covers API for requests, responses, cookies, sessions, headers; excludes advanced networking and async features. |
| statsmodels | StatModeler | Statistical Modeling Library | Provides econometric and statistical modeling tools, including regression, time series, and hypothesis testing. | Focuses on classical statistical analysis and diagnostics; excludes modern machine learning and deep learning. |

```

    ],
    "task_query": "You are testing an algorithm that applies migration operations to GIS models, ensuring
                  that spatial indexes on spatial fields are properly created, enabled, disabled, or removed as
                  dictated by the migration specifications.",
    "id": "django-0109"
}

```

Each task is represented by (i) its repository category and file location, (ii) the associated test class and functions, and (iii) a natural-language query summarizing the algorithm under test.

Given such a task, the benchmark provides the **algorithm description**, its corresponding **input–output ground truth**, and the **test method**. Evaluation is then conducted along two dimensions: (1) *Algorithm Presence* — whether the generated repository contains an implementation that matches the target algorithm, and (2) *Algorithm Correctness* — whether the adapted tests pass against the generated implementation, reflecting functional accuracy. This dual perspective allows us to measure both coverage of algorithmic functionality and the reliability of generated implementations.

D.3 Agent Pipeline

The evaluation employs a three-stage agent pipeline to connect task descriptions with generated repositories and derive executable judgments of success.

Stage 1: Localization. Given a task and its algorithmic description, the agent first explores the generated repository to locate candidate functions or classes that may implement the target algorithm. This step uses the exploration tools detailed in Appendix C.1, and produces a set of potentially relevant code anchors.

Stage 2: Majority-Vote Validation. To verify whether the localized candidates truly correspond to the target algorithm, we employ a majority-voting mechanism with a large language model (LLM). Each candidate is evaluated five times; the majority outcome is taken as the decision. If the validation fails, the pipeline triggers a re-localization

attempt. The localization–validation loop is retried up to three times; if all attempts fail, the repository is judged to lack an implementation of the algorithm.

Stage 3: Test Adaptation and Execution. For validated candidates, the agent then adapts the task’s reference test code. Concretely, the provided ground-truth test (including inputs, outputs, and checking methods) is rewritten to match the naming and structural conventions of the localized function or class. The adapted test is executed, and its outcome determines whether the generated implementation is functionally correct.

This pipeline ensures that evaluation captures both *coverage* (whether an algorithm is present in the generated repository) and *correctness* (whether its implementation passes the adapted tests).

D.3.1 Metrics

To comprehensively evaluate the generated repositories, we adopt a multi-dimensional set of metrics that capture four complementary aspects: *functionality alignment*, *novelty*, *execution accuracy*, and *code scale*. The motivation is to move beyond a single success/failure judgment and instead characterize (i) whether the right algorithms are generated, (ii) whether new functionalities are introduced, (iii) whether these implementations actually work, and (iv) at what level of scale and complexity they are realized. Together, these metrics provide a holistic view of the strengths and limitations of different models.

Functionality Coverage. The first question is whether a model can reproduce the expected range of functionalities in a target repository. We extract feature descriptions from both ground-truth repositories and generated repositories, and define a reference set of categories $\mathcal{C} = \{c_1, \dots, c_K\}$ based on official documentation and developer guidelines. Generated functionalities $\mathcal{G} = \{g_1, \dots, g_N\}$ are obtained either from structured intermediate outputs (for agent-based methods) or directly from raw code (for baseline models). To align generated features with reference categories, we perform K-Means clustering with \mathcal{C} as fixed centroids, plus an additional centroid c_{OOD} for out-of-distribution features. Each generated feature g_i is mapped to $f(g_i) \in \mathcal{C} \cup \{c_{\text{OOD}}\}$, with assignments further refined by an LLM-as-Judge to reduce semantic drift. Coverage is then defined as the fraction of reference categories that are “hit” by at least one generated feature:

$$\text{Coverage} = \frac{1}{|\mathcal{C}|} \sum_{j=1}^K \mathbb{1} [\exists g_i \in \mathcal{G}, f(g_i) = c_j]. \quad (1)$$

This metric quantifies how well the generated repository aligns with the intended functionality footprint.

Functionality Novelty. Coverage alone cannot distinguish between a model that simply memorizes existing categories and one that proposes extensions. To capture creativity and diversity, we measure the proportion of generated functionalities that fall outside the reference taxonomy. Specifically, novelty is the fraction of generated nodes assigned to the out-of-distribution centroid c_{OOD} :

$$\text{Novelty} = \frac{1}{|\mathcal{G}|} \sum_{i=1}^N \mathbb{1} [f(g_i) = c_{\text{OOD}}]. \quad (2)$$

High novelty indicates a tendency to introduce new capabilities, though such capabilities may or may not be useful. This metric is therefore best interpreted jointly with accuracy (below).

Functionality Accuracy. Even if a repository covers the right categories, the implementations must be correct. We therefore evaluate repository-specific tasks by checking whether generated code passes adapted test cases. Two complementary statistics are reported:

- **Voting Rate** — the fraction of tasks where the localization–validation pipeline confirms that an implementation of the target algorithm is present. This measures algorithm *presence*.
- **Success Rate** — the fraction of tasks where the adapted tests execute successfully. This measures algorithm *correctness*.

Together, these metrics disentangle whether errors stem from missing functionality versus incorrect implementation.

Code-Level Statistics. Finally, we report statistics on the scale and complexity of generated codebases. This helps distinguish minimal solutions from more realistic, full-fledged repositories. We compute these metrics over filtered Python source files, excluding directories unrelated to core functionality (e.g., `tests`, `examples`, `benchmarks`). The reported quantities are:

- **File Count:** number of valid source files, reflecting modular spread;
- **Normalized LOC:** effective lines of code after removing comments, docstrings, and blank lines, capturing implementation size;
- **Code Token Count:** number of tokens in normalized code, measured with a standard tokenizer, reflecting lexical complexity.

By jointly considering these four dimensions (coverage, novelty, accuracy in terms of presence and correctness, and scale), we obtain a nuanced evaluation of generated repositories. This design ensures that models are rewarded not only for producing functional code, but also for producing diverse, accurate, and realistically sized repositories.

D.4 Ground-Truth Taxonomy for Coverage and Novelty Calculation

```
{
  "Supervised learning": {
    "Linear Models": [],
    "Kernel ridge regression": [],
    "Support Vector Machines": {
      "SVC": [],
      "SVR": [],
      "OneClassSVM": []
    },
    "Nearest Neighbors": [],
    "Gaussian Processes": [],
    "Cross decomposition": [],
    "Naive Bayes": [],
    "Decision Trees": [],
    "Ensembles": [],
    "Multiclass and multioutput algorithms": [],
    "Feature selection": [],
    "Semi-supervised learning": [],
    "Isotonic regression": [],
    "Probability calibration": [],
    "Neural network models (supervised)": []
  },
  "Unsupervised learning": {
    "Gaussian mixture models": [],
    "Manifold learning": [],
    "Clustering": [],
    "Biclustering": [],
    "matrix factorization problems": [],
    "Covariance estimation": [],
    "Novelty and Outlier Detection": [],
    "Density Estimation": [],
    "Neural network models (unsupervised)": []
  },
  "Model selection and evaluation": {
    "Cross-validation": [],
    "Tuning the hyper-parameters of an estimator": [],
    "Validation curves": {
      "Classification Metrics": [],
      "Regression Metrics": [],
      "Clustering Metrics": []
    }
  },
  "Inspection": {
    "Permutation feature importance": []
  },
  "Dataset transformations": {
    "Feature extraction": {
      "Feature hashing": [],
      "Text feature extraction": [],

```

```

    "Image feature extraction": []
  },
  "Preprocessing data": {
    "Scaling and Normalization": [],
    "Discretization and Binarization": [],
    "Polynomial and Non-linear Feature Engineering": [],
    "Categorical Feature Encoding": [],
    "Missing Value Imputation": [],
    "Kernel and Matrix Centering": []
  },
  "Unsupervised dimensionality reduction": [],
  "Random Projection": [],
  "Kernel Approximation": [],
  "Pairwise metrics, Affinities and Kernels": []
},
"Dataset loading utilities": {},
"Model persistence": []
}

```

```

{
  "requests": {
    "Core Request Features": {
      "HTTP Method Support": [],
      "URL and Query Handling": [],
      "Request Body Construction": [],
      "Custom Headers": []
    },
    "Response Handling": {
      "Response Body Access": [],
      "JSON Processing": [],
      "Response Metadata": [],
      "Cookie Handling": []
    },
    "Session Management": {
      "Session Persistence": [],
      "Session Customization": []
    },
    "Advanced Configuration": {
      "Timeouts and Retries": [],
      "Redirect Control": [],
      "Streaming and Chunking": [],
      "Authentication Support": [],
      "Event Hooks": []
    },
    "Security and Transport": {
      "SSL Verification": [],
      "Client Certificates": [],
      "Transport Control": [],
      "Proxy Support": []
    },
    "Compliance and Encoding": {
      "Encoding Handling": [],
      "Standards Compliance": [],
      "Blocking Behavior": []
    }
  }
}

```


E Experiment Results

E.1 Baseline Configurations

To ensure fair comparison, we evaluate three representative systems for repository synthesis: **MetaGPT**, **ChatDev**, and **Paper2Code**, together with several single-agent LLM baselines. All methods are run with their official or default configurations.

MetaGPT. MetaGPT is a multi-agent framework that simulates a software company by assigning roles such as Product Manager, Architect, Project Manager, Engineer, and Tester. The agents collaborate following predefined Standard Operating Procedures to complete planning, design, implementation, and debugging.

ChatDev. ChatDev also follows a company-style organization, where agents take charge of requirement analysis, coding, testing, and review. It uses a chat-based interaction mechanism to coordinate between stages. We run ChatDev with its default settings.

Paper2Code. Paper2Code is a fixed workflow system designed to convert machine learning papers into executable repositories. It follows a three-stage pipeline of planning, analysis, and generation, which we execute sequentially using the default setup.

Vibe-Coding Agent (OpenHands, Codex, Claude Code, Gemini CLI). For comparison with standalone LLM systems, we configure each model with a maximum of 30 iterations. The first round is initialized with the repository description. In each subsequent round, the model receives a fixed self-reflection prompt:

Please check whether the current repository still has any features that could be enhanced or any missing functionality that needs to be added. If there are no further improvements, or if you consider the task complete, please reply with "yes" only. If there are still potential enhancements or improvements to be made, please continue working on them, and do not reply with "yes" just because you are concerned about complexity.

E.2 Detailed Experiment Results

We report the results of different methods on six repositories. For each repository, the methods are evaluated under the same settings to enable direct comparison.

Table 11: Performance on the **MLKit-Py** "Nov." denotes the novelty rate; the number in parentheses is Novel/Total, where Novel is the number of novel functionalities and Total is the total number of planned functionalities.

| Agent | Model | Cov. (%) ↑ | Nov. (%) (Novel/Total) ↑ | Pass. / Vot. (%) ↑ | Files ↑ | LOC ↑ | Tokens ↑ |
|-----------------|------------------|-------------|--------------------------|--------------------|---------|-------|----------|
| MetaGPT | o3-mini | 14.9 | 0.0 (0.0/13.0) | 6.3 / 7.3 | 3.0 | 95.0 | 928.0 |
| | Qwen3-Coder | 19.2 | 0.0 (0.0/23.0) | 9.9 / 12.0 | 8.0 | 170.0 | 1718 |
| ChatDev | o3-mini | 8.5 | 14.3 (2/14) | 6.3 / 7.3 | 6 | 163 | 2064 |
| | Qwen3-Coder | 12.8 | 0.0 (0/49) | 10.5 / 11.5 | 7 | 280 | 3100 |
| OpenHands | o3-mini | 31.9 | 0.0 (0/39) | 11.5 / 13.6 | 14 | 272 | 2499 |
| | Qwen3-Coder | 34.0 | 0.0 (0/48) | 11.0 / 14.0 | 26 | 1020 | 10213 |
| Paper2Code | o3-mini | 25.5 | 0.0 (0/41) | 17.8 / 19.9 | 5 | 564 | 6346 |
| | Qwen3-Coder | 31.9 | 0.0 (0/118) | 18.8 / 24.6 | 12 | 1710 | 20821 |
| Codex CLI | o3 pro | 31.9 | 0.0 (0/59) | 11.0 / 16.9 | 14 | 829 | 8344 |
| Gemini CLI | gemini 2.5 pro | 59.6 | 0.0 (0/141) | 0.0 / 33.5 | 19 | 2316 | 24782 |
| Claude Code CLI | claude 4 sonnet | 59.6 | 0.0 (0/163) | 27.5 / 42.4 | 31 | 3559 | 37056 |
| Gold Projects | Human Developers | - | - | 85.1 / 98.3 | 185 | 65972 | 592187 |
| ZeroRepo | o3-mini | 97.9 | 4.7 (54/1258) | 73.5 / 78.7 | 266 | 31596 | 351554 |
| | Qwen3-Coder | 85.1 | 15.0 (176/1233) | 63.6 / 74.6 | 642 | 60553 | 741634 |

E.3 Examples of Coverage Calculation and Novelty Assessment

In this subsection, we provide examples of how coverage and novelty are computed from the constructed RPG, illustrating category alignment for coverage and out-of-distribution detection for novelty.

Table 12: Performance on the **HttpEasy** repo. "Nov." denotes the novelty rate; the number in parentheses is Novel/Total, where Novel is the number of novel functionalities and Total is the total number of planned functionalities.

| Agent | Model | Cov. (%) ↑ | Nov. (%) (Novel/Total) ↑ | Pass. / Vot. (%) ↑ | Files ↑ | LOC ↑ | Tokens ↑ |
|-----------------|------------------|--------------|--------------------------|--------------------|---------|-------|----------|
| MetaGPT | o3-mini | 22.7 | 0.0 (0/12) | 5.0 / 15.0 | 1 | 167 | 1802 |
| | Qwen3-Coder | 31.8 | 0.0 (0/17) | 20.0 / 25.0 | 4 | 175 | 2023 |
| ChatDev | o3-mini | 36.4 | 18.2 (2/11) | 15.0 / 15.0 | 3 | 177 | 2055 |
| | Qwen3-Coder | 40.9 | 3.5 (1/31) | 20.0 / 30.0 | 2 | 323 | 3151 |
| OpenHands | o3-mini | 22.7 | 0.0 (0/5) | 20.5 / 28.2 | 3 | 72 | 669 |
| | Qwen3-Coder | 31.8 | 0.0 (0/20) | 20.0 / 30.0 | 2 | 214 | 1960 |
| Paper2Code | o3-mini | 27.3 | 0.0 (0/18) | 0.0 / 24.2 | 5 | 192 | 1856 |
| | Qwen3-Coder | 50.0 | 2.7 (1/39) | 0.0 / 45.5 | 5 | 377 | 3965 |
| Codex CLI | o3 pro | 45.5 | 0.0 (0/19) | 14.0 / 28.0 | 1 | 197 | 1879 |
| Gemini CLI | gemini 2.5 pro | 59.1 | 3.1 (1/33) | 40.0 / 56.0 | 1 | 420 | 5407 |
| Claude Code CLI | claude 4 sonnet | 50.0 | 0.0 (0/21) | 36.0 / 42.0 | 2 | 436 | 4931 |
| Gold Projects | Human Developers | - | - | 72.3 / 87.2 | 17 | 2793 | 22297 |
| ZeroRepo | o3-mini | 100.0 | 2.05 (7/433) | 64.0 / 72.0 | 109 | 6192 | 61922 |
| | Qwen3-Coder | 95.5 | 0.3 (2/854) | 54.0 / 64.0 | 245 | 15559 | 165051 |

Table 13: Performance on the **PyWebEngine** repo. "Nov." denotes the novelty rate; the number in parentheses is Novel/Total, where Novel is the number of novel functionalities and Total is the total number of planned functionalities.

| Agent | Model | Cov. (%) ↑ | Nov. (%) (Novel/Total) ↑ | Pass. / Vot. (%) ↑ | Files ↑ | LOC ↑ | Tokens ↑ |
|-----------------|------------------|-------------|--------------------------|--------------------|---------|--------|----------|
| MetaGPT | o3-mini | 27.1 | 0.0 (0/52) | 0.0 / 13.5 | 2 | 421 | 3733 |
| | Qwen3-Coder | 18.8 | 0.0 (0/52) | 0.0 / 9.2 | 9 | 238 | 1928 |
| ChatDev | o3-mini | 25.0 | 0.0 (0/40) | 0.0 / 14.2 | 8 | 372 | 3185 |
| | Qwen3-Coder | 27.1 | 0.0 (0/49) | 0.0 / 12.1 | 11 | 679 | 5950 |
| OpenHands | o3-mini | 31.3 | 2.0 (1/55) | 0.0 / 14.2 | 18 | 304 | 2628 |
| | Qwen3-Coder | 25.0 | 0.0 (0/52) | 0.0 / 19.1 | 13 | 427 | 3996 |
| Paper2Code | o3-mini | 27.1 | 0.0 (0/46) | 0.0 / 15.6 | 11 | 619 | 6342 |
| | Qwen3-Coder | 43.8 | 0.0 (0/103) | 0.0 / 19.9 | 10 | 1761 | 16076 |
| Codex CLI | o3 pro | 39.6 | 0.0 (0/88) | 12.1 / 26.7 | 2 | 769 | 7751 |
| Gemini CLI | gemini 2.5 pro | 45.8 | 0.3 (1/318) | 7.6 / 48.1 | 45 | 2975 | 27655 |
| Claude Code CLI | claude 4 sonnet | 64.6 | 38.1 (669/2165) | 33.9 / 66.1 | 80 | 34302 | 317883 |
| Gold Projects | Human Developers | - | - | 81.6 / 86.5 | 681 | 109457 | 917622 |
| ZeroRepo | o3-mini | 79.2 | 38.2 (566/1680) | 74.1 / 84.4 | 430 | 27647 | 275782 |
| | Qwen3-Coder | 68.8 | 18.1 (244/1561) | 56.4 / 64.8 | 521 | 48058 | 539052 |

Analysis of Coverage Examples. These examples demonstrate that our coverage metric provides a reasonable allocation of generated functionalities to reference categories. Core areas such as regression, classification, clustering, and preprocessing are consistently captured, while supporting utilities (e.g., normalization, imputation) are distributed into their respective modules without overlap or misplacement. This validates the soundness of our metric design for assessing functional completeness. Moreover, the RPG ensures that functionalities are not only well aligned with reference categories but also diversified across them, highlighting its effectiveness as a planning substrate for repository-level generation.

```
{
  "SVR": [
    "NuSVR"
  ],
  "Gaussian mixture models": [
    "gmm expectation maximization",
    "dp gaussian mixture"
  ],
  "Scaling and Normalization": [
    "quantile scaling",
```

Table 14: Performance on the **TableKit** repo. "Nov." denotes the novelty rate; the number in parentheses is Novel/Total, where Novel is the number of novel functionalities and Total is the total number of planned functionalities.

| Agent | Model | Cov. (%) ↑ | Nov. (%) (Novel/Total) ↑ | Pass. / Vot. (%) ↑ | Files ↑ | LOC ↑ | Tokens ↑ |
|-----------------|------------------|-------------|--------------------------|--------------------|---------|--------|----------|
| MetaGPT | o3-mini | 13.2 | 0.0 (0/21) | 0.0 / 11.5 | 1 | 186 | 1814 |
| | Qwen3-Coder | 6.6 | 0.0 (0/17) | 0.0 / 6.4 | 3 | 133 | 1453 |
| ChatDev | o3-mini | 21.1 | 0.0 (0/36) | 0.0 / 15.0 | 2 | 332 | 3517 |
| | Qwen3-Coder | 19.7 | 0.0 (0/54) | 0.0 / 0.0 | 6 | 918 | 9168 |
| OpenHands | o3-mini | 11.8 | 0.0 (0/26) | 0.0 / 18.1 | 6 | 193 | 1753 |
| | Qwen3-Coder | 11.8 | 0.0 (0/23) | 0.0 / 12.1 | 2 | 174 | 1914 |
| Paper2Code | o3-mini | 17.1 | 9.4 (5/53) | 0.0 / 23.5 | 7 | 529 | 5325 |
| | Qwen3-Coder | 17.1 | 0.0 (0/61) | 6.2 / 20.4 | 9 | 1886 | 19337 |
| Codex CLI | o3 pro | 11.8 | 0.0 (0/23) | 21.1 / 30.9 | 2 | 552 | 6299 |
| Gemini CLI | gemini 2.5 pro | 44.7 | 0.0 (0/117) | 38.6 / 48.5 | 15 | 1249 | 12242 |
| Claude Code CLI | claude 4 sonnet | 52.6 | 0.0 (0/191) | 53.1 / 77.7 | 11 | 8509 | 83834 |
| Gold Projects | Human Developers | - | - | 90.6 / 94.0 | 217 | 106447 | 943873 |
| ZeroRepo | o3-mini | 72.4 | 21.1 (306/1701) | 81.4 / 88.3 | 477 | 37331 | 395536 |
| | Qwen3-Coder | 65.8 | 13.9 (178/1500) | 48.0 / 64.8 | 347 | 32387 | 389886 |

Table 15: Performance on the **StatModeler** repo. "Nov." denotes the novelty rate; the number in parentheses is Novel/Total, where Novel is the number of novel functionalities and Total is the total number of planned functionalities.

| Agent | Model | Cov. (%) ↑ | Nov. (%) (Novel/Total) ↑ | Pass. / Vot. (%) ↑ | Files ↑ | LOC ↑ | Tokens ↑ |
|-----------------|------------------|-------------|--------------------------|--------------------|---------|-------|----------|
| MetaGPT | o3-mini | 11.4 | 0.0 (0/19) | 5.6 / 6.1 | 6 | 228 | 2330 |
| | Qwen3-Coder | 5.7 | 0.0 (0/10) | 0.0 / 2.8 | 13 | 437 | 5435 |
| ChatDev | o3-mini | 10.2 | 21.1 (8/38) | 1.1 / 9.5 | 9 | 726 | 9644 |
| | Qwen3-Coder | 11.4 | 0.0 (0/18) | 3.2 / 7.7 | 6 | 320 | 3797 |
| OpenHands | o3-mini | 13.6 | 0.0 (0/32) | 7.9 / 9.0 | 9 | 335 | 3338 |
| | Qwen3-Coder | 14.8 | 0.0 (0/27) | 9.5 / 12.8 | 5 | 670 | 8476 |
| Paper2Code | o3-mini | 12.5 | 21.6 (8/29) | 0.0 / 10.7 | 9 | 813 | 9793 |
| | Qwen3-Coder | 13.6 | 30.0 (12/50) | 3.2 / 14.0 | 8 | 1179 | 13519 |
| Codex CLI | o3 pro | 20.5 | 0.0 (0/23) | 8.2 / 9.9 | 9 | 709 | 8473 |
| Gemini CLI | gemini 2.5 pro | 23.7 | 0.0 (0/55) | 13.5 / 23.2 | 6 | 736 | 8063 |
| Claude Code CLI | claude 4 sonnet | 34.1 | 0.0 (0/191) | 18.4 / 27.8 | 28 | 4043 | 46182 |
| Gold Projects | Human Developers | - | - | 87.2 / 96.2 | 271 | 83325 | 893824 |
| ZeroRepo | o3-mini | 77.3 | 15.6 (143/1021) | 76.4 / 81.1 | 220 | 24141 | 294292 |
| | Qwen3-Coder | 77.3 | 8.2 (83/1113) | 66.2 / 73.9 | 436 | 47370 | 598058 |

```

    "scale to [0, 1]",
    "z-score scaling",
    "IQR scaling"
  ],
  "Missing Value Imputation": [
    "mean imputation",
    "matrix completion imputation",
    "impute using K-nearest neighbors",
    "impute with global median",
    "impute missing data",
  ],
  "Ensembles": [
    "light gradient boosting",
    "HistGradientBoosting",
    "bagging classification trees",
    "random forest",
    "LightGBM",
    "CatBoost",
    "XGBoost"
  ]

```

```

],
"Clustering Metrics": [
  "density peak clustering",
  "gap statistic",
  "silhouette score calculation",
  "inertia calculation"
],
"Naive Bayes": [
  "multinomial naive bayes",
  "bernoulli naive bayes",
  "gaussian naive bayes"
],
"Linear Models": [
  "ridge regression",
  "lasso regression",
  "huber regression",
  "ransac regression"
],
"SVC": [
  "soft margin SVM",
  "hard margin SVM",
  "SVM with precomputed kernel"
]
}

```

```

{
  "Proxy Support": [
    "rotate proxy list",
    "auto detect system proxy",
    "custom dns resolver integration"
  ],
  "HTTP Method Support": [
    "send POST request",
    "GET request with cookies",
    "send DELETE request",
    "PUT with JSON payload"
  ],
  "URL and Query Handling": [
    "encode path segments",
    "parse query string",
    "normalize request url"
  ],
  "Redirect Control": [
    "auto follow redirects",
    "limit redirect chain"
  ],
  "Authentication Support": [
    "send basic auth",
    "include oauth2 bearer token",
    "refresh auth token"
  ],
  "Timeouts and Retries": [
    "set request timeout",
    "apply exponential backoff",
    "custom retry hook"
  ],
  "JSON Processing": [
    "auto deserialize json",
    "validate json schema",
    "serialize dict to JSON"
  ],
}

```

```

"SSL Verification": [
  "ssl hostname verification",
  "load custom certificates"
],
"Streaming and Chunking": [
  "process chunked response",
  "resume file download support"
]
}

```

Analysis of Novelty Examples The novelty cases illustrate two key observations. First, novelty captures meaningful extensions rather than random noise: in MLKit-Py, we see coherent additions such as *Prophet forecasting*, *STL decomposition*, and *genetic programming feature synthesis*, while in StatModeler new capabilities include *vector autoregression* and *Cox proportional hazards models*. Second, the new functionalities proposed by the RPG remain reasonable within the target domain: they extend statistical modeling, optimization, or robustness analysis in ways that align with real-world software evolution. Together, these examples confirm that the RPG supports not only stable replication of reference repositories but also the introduction of coherent and domain-consistent innovations.

```

{
  "new_features": [
    "vector autoregression model",
    "forecasting with Prophet",
    "genetic programming feature synthesis",
    "multi-objective bayesian optimization",
    "online learning",
    "apriori association mining",
    "Cox proportional hazards model",
    "STL decomposition",
    "temporal drift detection",
    "fuzz testing",
    "interactive dashboards",
    "NoSQL queries",
    "pareto optimization",
    "demographic parity test",
    "secure argument parsing",
    ...
  ]
}

```

```

{
  "new_features": [
    "vector autoregression model",
    "forecasting with Prophet",
    "genetic programming feature synthesis",
    "multi-objective bayesian optimization",
    "online learning",
    "apriori association mining",
    "Cox proportional hazards model",
    "STL decomposition",
    "temporal drift detection",
    "fuzz testing",
    "interactive dashboards",
    "NoSQL queries",
    "pareto optimization",
    "demographic parity test",
    "secure argument parsing",
    ...
  ]
}

```

```
]
}
```

E.4 Examples of Localization Behavior

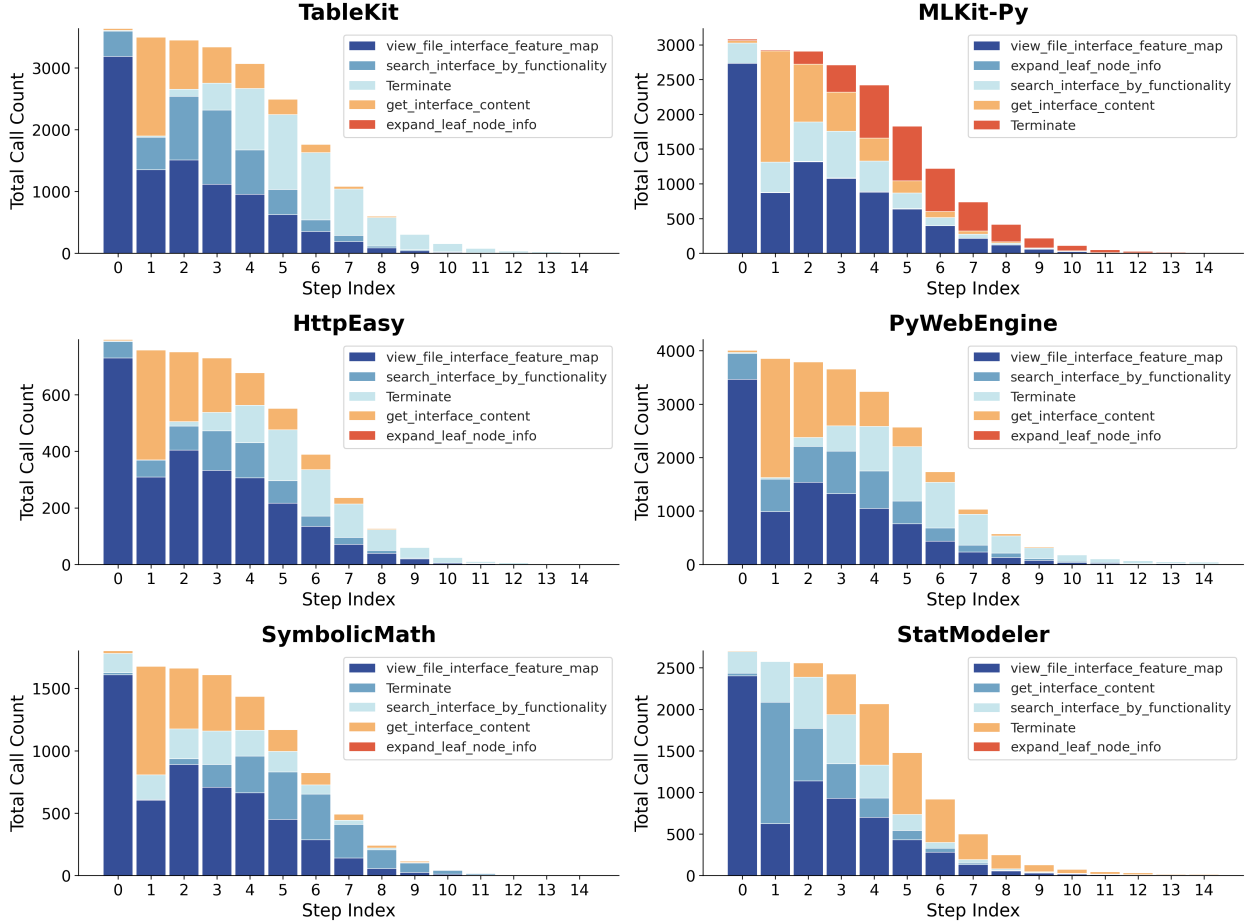


Figure 13: Aggregated function call frequency distribution across localization steps in all repositories using o3-mini.

Graph guidance structures localization into systematic search. Figure 13 shows that with graph guidance, localization behavior follows a structured CCG pattern (Coarse Search → Content Inspection → Global Graph Exploration). The agent begins by traversing the RPG at a coarse level to identify high-level candidates, then inspects content-rich nodes for detailed signals, and finally explores semantically related structures across the graph. Termination calls rise as the search converges. This progression indicates that the RPG reshapes the agent’s behavior into a systematic, relation-aware search process, replacing ad hoc or repetitive probing.