

Is Reasoning All You Need? An Empirical Study of Retrieval and Structured Reasoning for Python Class-Level Docstring Generation

Anonymous ACL submission

Abstract

Automated docstring generation is a high-fidelity software engineering task where factual correctness, interface coverage, and efficiency are critical. Recent work increasingly applies generation-time reasoning strategies—such as Chain-of-Thought (CoT), Tree-of-Thought (ToT), and Graph-of-Thought (GoT)—yet their benefits relative to retrieval-based grounding for code documentation remain unclear. We present a controlled empirical evaluation of 12 Python class-level docstring generation strategies, combining three architectural families—Plain LLM, Retrieval-Augmented Generation (RAG), and Iterative Critique RAG—with four reasoning modes: Base, CoT, ToT, and GoT. Strategies are evaluated using lexical and semantic metrics (ROUGE, BLEU, BERTScore), code-specific coverage measures, an LLM-as-a-Judge faithfulness metric, and latency-based cost analysis.

Our results reveal three consistent findings. First, retrieval is the primary driver of factual reliability: a simple RAG baseline achieves 73% faithfulness, compared to 50% for Plain LLMs and 57% for Iterative Critique RAG (+46% relative improvement). Second, reasoning-heavy strategies show diminishing returns: ToT and GoT increase latency by up to 15× without statistically significant gains in faithfulness or semantic similarity ($\leq +2\text{-}3\%$ BERTScore). Third, SimpleRAG offers the best cost-quality trade-off, generating high-quality docstrings in 6 seconds versus 90 seconds for Tree-of-Thought pipelines. Overall, our findings indicate that for Python class-level docstring generation—a static, context-bound task—efficient retrieval-based grounding is more effective than generation-time reasoning. We release a reproducible, cost-aware benchmark to support principled evaluation of automated documentation systems.

1 Introduction

Large Language Models (LLMs) have enabled substantial progress in software engineering tasks such as code generation, explanation, and summarization. To improve reliability, recent work increasingly augments LLMs with external context, explicit reasoning, and self-correction mechanisms. Among these approaches, Retrieval-Augmented Generation (RAG) has become a standard technique for grounding model outputs in external documentation rather than relying solely on parametric knowledge.

Automatic docstring generation is a particularly important application of LLMs. Unlike open-ended summarization, docstring generation is a high-fidelity task in which generated descriptions must accurately reflect code behavior, parameters, and return values. Hallucinated or incorrect docstrings can mislead developers and propagate errors into downstream tooling, making this task a stringent benchmark for evaluating factual grounding in LLM-based systems.

While modern LLMs can produce fluent docstrings without external context, standalone generation frequently yields generic or incorrect descriptions, especially for unfamiliar or domain-specific code (Zhang et al., 2025b). Prior work therefore adopts RAG-based approaches that ground generation in style guides, API documentation, and coding conventions, demonstrating improved factual correctness through retrieval-based grounding. Figure 5 illustrates the qualitative difference between docstrings generated with and without external context.

Concurrently, recent research emphasizes generation-time reasoning strategies—such as Chain-of-Thought (CoT), Tree-of-Thought (ToT), Graph-of-Thought (GoT), and agentic self-correction—as general solutions to hallucination. These techniques have demonstrated strong per-

083 performance on tasks requiring multi-step reasoning
084 and planning and are increasingly being applied to
085 code documentation pipelines. This trend raises
086 a critical but under-examined question for both
087 researchers and practitioners: **whether the sub-**
088 **stantial computational and latency overhead in-**
089 **troduced by explicit reasoning is justified for**
090 **documentation tasks that are structurally sim-**
091 **ple, context-bound, and frequently deployed in**
092 **latency-sensitive developer tooling.**

093 Despite the widespread adoption of both RAG
094 and reasoning-based prompting, their interaction
095 in structured, context-bound documentation tasks
096 remains poorly understood. Existing studies on
097 RAG for code documentation (Zhang et al., 2025c)
098 focus on limited retrieval variants, while reasoning-
099 centric work largely targets open-ended reasoning
100 tasks rather than extractive documentation.

101 In this work, we present a systematic empiri-
102 cal benchmark of retrieval and reasoning strate-
103 gies for Python class-level docstring generation.
104 We evaluate twelve configurations formed by com-
105 bining three architectural families—Plain LLMs,
106 standard RAG, and iterative self-correction—with
107 four reasoning modalities: Base generation, CoT,
108 ToT, and GoT. This factorial design enables con-
109 trolled isolation of the effects of retrieval, explicit
110 reasoning, and self-verification. We conduct exper-
111 iments on a curated dataset of Python class doc-
112 strings and evaluate performance using lexical and
113 semantic metrics: ROUGE (Grusky, 2023), BLEU
114 (Papineni et al., 2002), and BERTScore (Zhang*
115 et al., 2020); code-specific coverage measures (pa-
116 rameter, return, and exception coverage); and a
117 reference-free LLM-as-a-Judge (Bavaresco et al.,
118 2025) faithfulness metric to detect hallucinations
119 and unsupported claims. We additionally report in-
120 ference latency to analyze cost–quality trade-offs.

121 Our results show that retrieval-based grounding
122 is the dominant driver of reliable docstring gen-
123 eration: a lightweight RAG baseline achieves a
124 Faithfulness score of 73%, substantially outper-
125 forming Plain LLMs (50%) and Self-Correction
126 RAG (57%), corresponding to a 46% relative im-
127 provement over non-retrieval generation. In con-
128 trast, reasoning-heavy strategies such as Tree-of-
129 Thought and Graph-of-Thought exhibit diminish-
130 ing returns, increasing inference latency by up to
131 15× without statistically significant gains in faithful-
132 ness or semantic accuracy ($\leq +2\text{-}3\%$ BERTScore).
133 From an efficiency perspective, SimpleRAG offers
134 the most favorable quality–cost trade-off, gener-

ating high-quality docstrings in approximately 6
seconds per sample compared to roughly 90 sec-
onds for Tree-of-Thought pipelines—a 93% re-
duction in latency. Collectively, these findings
demonstrate that for code documentation, an in-
herently static and context-bound task, efficient
retrieval-based grounding is substantially more ef-
fective than generation-time reasoning.

2 Related Work 143

2.1 Automated Code Documentation 144

145 Automated code documentation, including sum-
146 maries and docstrings, has been widely studied to
147 improve software maintainability and developer
148 productivity. Neural approaches to code summa-
149 rization demonstrate that models trained jointly on
150 code and natural language can generate fluent de-
151 scriptions of program behavior. For example, Liu et
152 al. (Liu et al., 2021) propose a retrieval-augmented
153 summarization mechanism that integrates code
154 semantics with dense representations, reflecting
155 the evolution of retrieval-based techniques beyond
156 simple prompt tuning. However, recent studies
157 show that even strong code-specialized models fre-
158 quently hallucinate undocumented parameters or
159 behaviors when relying solely on parametric knowl-
160 edge (Dainese et al., 2024; Sundaram et al., 2025).
161 These findings motivate architectures that explicitly
162 ground generation in external context.

2.2 Retrieval-Augmented Generation for Code Tasks 163

164 Retrieval-Augmented Generation (RAG) improves
165 factual grounding by injecting external, non-
166 parametric information into LLM prompts at infer-
167 ence time. RAG has been studied in both general
168 NLP summarization and code-related tasks. For in-
169 stance, (Zhang et al., 2025a) introduces the Graph
170 of Records, which enriches RAG using historical
171 LLM responses for long-context summarization,
172 suggesting that richer context representations can
173 benefit global documentation tasks. Several exten-
174 sions to basic RAG—such as corrective retrieval,
175 query reformulation, and source fusion—have been
176 proposed. However, evaluations often focus on
177 single architectural variants or rely primarily on
178 surface-level text similarity metrics, leaving open
179 questions about reliability, efficiency, and faithful-
180 ness in structured documentation tasks. 181

| | | |
|-----|--|-----|
| 182 | 2.3 Advanced Reasoning Strategies for LLMs | 231 |
| 183 | Beyond retrieval, explicit generation-time reason- | 232 |
| 184 | ing strategies have been introduced to improve | 233 |
| 185 | LLM performance on complex tasks. Chain-of- | 234 |
| 186 | Thought (CoT) prompting (Wei et al., 2022) en- | 235 |
| 187 | courages step-by-step reasoning, while Tree-of- | 236 |
| 188 | Thought (ToT) (Yao et al., 2023) and Graph-of- | 237 |
| 189 | Thought (GoT) (Besta et al., 2024) explore and | 238 |
| 190 | aggregate multiple reasoning paths. These meth- | 239 |
| 191 | ods show strong gains on tasks requiring multi- | 240 |
| 192 | step inference or planning. However, their ef- | |
| 193 | fectiveness for context-bound, extraction-oriented | |
| 194 | tasks—such as code summarization and documenta- | |
| 195 | tion—remains unclear. When relevant context is | |
| 196 | already supplied through retrieval, additional rea- | |
| 197 | soning may introduce speculative inference rather | |
| 198 | than improve factual grounding. | |
| 199 | 2.4 Evaluation of Code Summarization and | |
| 200 | LLM Outputs | |
| 201 | Standard metrics such as BLEU and ROUGE are | |
| 202 | widely used to evaluate generated summaries and | |
| 203 | documentation but often correlate poorly with hu- | |
| 204 | man judgments. We discuss this line of work in | |
| 205 | two parts: first, the limitations of surface-level | |
| 206 | automatic metrics for code documentation, and | |
| 207 | second, recent advances in LLM-based evalua- | |
| 208 | tion frameworks designed to better capture factual | |
| 209 | grounding and faithfulness. Recent work therefore | |
| 210 | emphasizes efficiency and cost alongside correct- | |
| 211 | ness. ShortenDoc (Yang et al., 2025), for example, | |
| 212 | demonstrates that compressing docstring-related | |
| 213 | prompts can significantly reduce inference costs | |
| 214 | without degrading performance. To better capture | |
| 215 | documentation quality, recent studies adopt code- | |
| 216 | specific metrics such as parameter coverage, re- | |
| 217 | turn coverage, and faithfulness scores that extend | |
| 218 | beyond token overlap. In parallel, the LLM-as- | |
| 219 | a-Judge paradigm has emerged as an alternative | |
| 220 | evaluation approach. Studies such as (Li et al., | |
| 221 | 2025a) show that LLM judges can assist in ranking | |
| 222 | and evaluating retrieval outputs, motivating their | |
| 223 | use for detecting subtle grounding errors missed by | |
| 224 | traditional metrics. | |
| 225 | 2.5 Research Gap | |
| 226 | While prior work has advanced retrieval- | |
| 227 | augmented generation, reasoning strategies, and | |
| 228 | evaluation methods, integrated comparisons of | |
| 229 | retrieval and reasoning for high-fidelity documen- | |
| 230 | tation remain limited. In particular, it is unclear | |
| | whether reasoning-enhanced pipelines improve | 231 |
| | faithfulness over retrieval-only approaches when | 232 |
| | factual grounding and efficiency are critical. More- | 233 |
| | over, most reasoning studies emphasize accuracy | 234 |
| | in isolation and do not evaluate hallucination | 235 |
| | behavior under retrieved context. Our study | 236 |
| | addresses this gap by systematically benchmarking | 237 |
| | retrieval and reasoning configurations for Python | 238 |
| | docstring generation across quality, faithfulness, | 239 |
| | and cost dimensions. | 240 |
| | 3 Methodology | 241 |
| | 3.1 Experimental Design | 242 |
| | We adopt a factorial experimental design that | 243 |
| | crosses three architectural families with four | 244 |
| | generation-time reasoning modes, yielding twelve | 245 |
| | docstring generation strategies. This design en- | 246 |
| | ables controlled analysis of (i) retrieval versus no | 247 |
| | retrieval, (ii) the marginal utility of explicit reason- | 248 |
| | ing, and (iii) the cost–benefit trade-offs introduced | 249 |
| | by iterative self-correction, while holding all other | 250 |
| | variables constant. All strategies are implemented | 251 |
| | within a unified, configuration-driven framework | 252 |
| | to ensure fair comparison and reproducibility. 1 | 253 |
| | shows the architecture diagram. | 254 |
| | 3.2 Task Definition | 255 |
| | The target task is Python class-level docstring gen- | 256 |
| | eration. Compared to function-level summariza- | 257 |
| | tion, class docstrings require synthesizing infor- | 258 |
| | mation across methods, attributes, and inheritance | 259 |
| | structures, making the task particularly sensitive to | 260 |
| | hallucination and omission. Each model receives | 261 |
| | the raw source code of a class and generates a PEP- | 262 |
| | style docstring describing its purpose, parameters, | 263 |
| | attributes, return behavior, and potential exceptions. | 264 |
| | This setting reflects realistic documentation work- | 265 |
| | flows and provides a stringent benchmark for evalu- | 266 |
| | ating factual grounding and interface completeness. | 267 |
| | 3.3 Architectural Families | 268 |
| | We evaluate three architectural paradigms that rep- | 269 |
| | resent increasing system complexity. | 270 |
| | Plain LLM serves as a no-retrieval baseline, re- | 271 |
| | lying solely on the model’s parametric knowledge | 272 |
| | and the input source code. This configuration estab- | 273 |
| | lishes a lower bound for documentation faithfulness | 274 |
| | in the absence of external grounding. | 275 |
| | Simple RAG follows the standard retrieval- | 276 |
| | augmented generation paradigm. A dense retriever | 277 |
| | selects the top- k relevant context chunks from an | 278 |

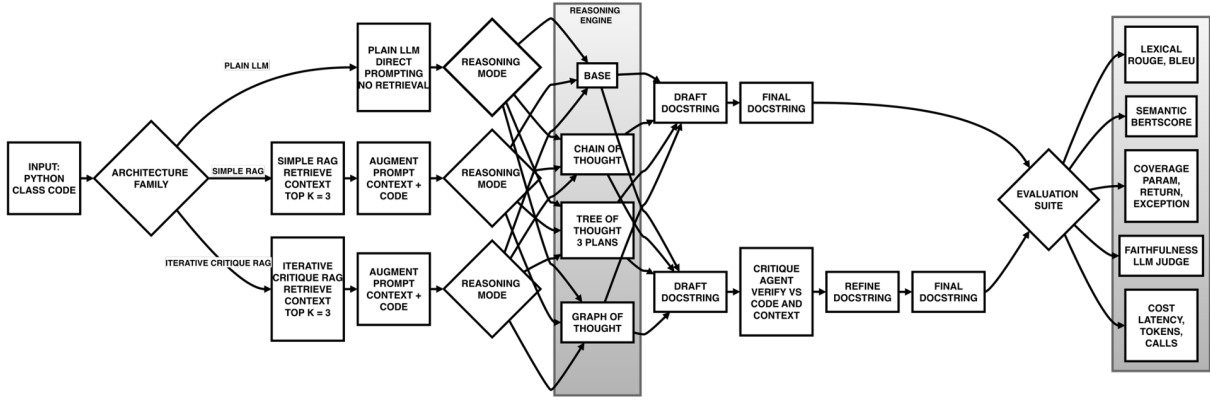


Figure 1: Experimental framework for Python docstring generation. Three architectural families combined with four reasoning modes yield twelve strategies, all evaluated using shared quality, coverage, faithfulness, and cost metrics.

indexed corpus and concatenates them with the source code prior to generation. This architecture isolates the effect of contextual grounding without introducing additional reasoning or control mechanisms.

Iterative Critique RAG applies a test-time generate–critique–refine loop in which the model revises an initial docstring based on self-identified omissions or unsupported claims. This heuristic self-correction approach, distinct from trained reflection-token frameworks such as SELF-RAG, incurs additional inference costs and is evaluated to determine whether lightweight self-review improves reliability beyond retrieval alone.

3.4 Reasoning Modes

Across all architectural families, we apply four prompting strategies (Li et al., 2025b) that vary the structure of generation-time reasoning.

The Base setting uses direct zero-shot (Izcard et al., 2022) prompting without explicit reasoning instructions. Chain-of-Thought (CoT) encourages step-by-step reasoning over parameters and behaviors. Tree-of-Thought (ToT) generates multiple candidate decomposition plans and synthesizes the most complete one. Graph-of-Thought (GoT) produces independent reasoning components (e.g., parameter analysis, return behavior, exception handling) and aggregates them into a final docstring. This axis enables systematic analysis of the benefits and computational costs of increasingly structured reasoning.

3.5 Evaluation Metrics

We employ a multi-dimensional evaluation framework covering documentation quality, completeness, faithfulness, and efficiency. Lexical similarity

is measured using ROUGE-1 and BLEU, while semantic alignment is assessed with BERTScore. Documentation completeness is evaluated using Parameter, Return, and Exception Coverage, which quantify whether relevant interface elements are explicitly documented.

To detect hallucinations not captured by similarity or coverage metrics, we introduce a reference-free Faithfulness Score using an independent LLM-as-a-Judge (Elmitwalli et al., 2025). The judge evaluates whether statements in the generated docstring are supported by the source code and, when applicable, retrieved context, producing scores from 0.0 (unsupported) to 1.0 (fully grounded). This score serves as our primary reliability metric. We additionally record inference latency, model call counts, and token usage to quantify computational cost and analyze quality–efficiency trade-offs.

3.6 Reproducibility

All experiments are driven by centralized configuration files that fix retrieval parameters, prompt templates, model settings, and evaluation criteria across strategies. This ensures that observed differences arise solely from architectural and reasoning variations rather than uncontrolled experimental factors.

4 Experimental Setup

4.1 Dataset Construction

We constructed a blind docstring generation benchmark from the rag-docstring dataset, consisting of 250 Python classes collected from diverse public GitHub repositories. For each class, we extract two aligned artifacts: (i) the raw class code and (ii) its original docstring. To enforce a blind gen-

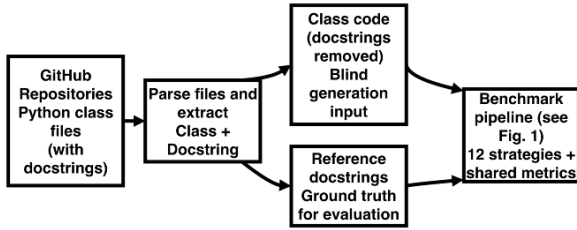


Figure 2: Dataset construction pipeline. Python classes are collected from GitHub, parsed to separate code and docstrings, and stripped of docstrings to create a blind generation setting; the extracted docstrings are used as reference annotations in the benchmarking framework (Figure 2).

eration setting, the docstring is removed from the input code prior to generation, while the extracted docstring is retained as the reference annotation. The resulting (code-without-docstring, reference-docstring) pairs form the evaluation dataset for all experiments. Figure 2 illustrates the dataset construction pipeline, and all samples are evaluated using the unified benchmarking framework shown in Figure 1.

4.2 Models and Hardware

All strategies are evaluated using a fixed set of models and hardware resources to avoid confounding effects.

- Generator LLM: Meta-Llama-3-8B-Instruct, quantized and served via Ollama, selected for its strong instruction-following capability and practical deployment footprint.
- Judge LLM: DeepSeek-Coder-V2, used exclusively as an independent evaluator for the LLM-as-a-Judge faithfulness metric.
- The embedding model all-MiniLM-L6-v2 (384 dimensions) was chosen to reflect realistic developer-tooling constraints rather than large, compute-heavy encoders.

All experiments were executed on a single NVIDIA A100 GPU, ensuring stable throughput and consistent latency measurements across all strategies. Reporting results on a high-performance accelerator allows observed cost differences to be attributed to architectural and reasoning choices rather than hardware limitations.

4.3 Retrieval Index and Parameters

To isolate the effects of reasoning and self-correction, retrieval parameters are fixed across all

retrieval-based strategies (Lewis et al., 2020). We use dense semantic retrieval with cosine similarity over embedded text chunks and retrieve the top- $k = 3$ documents. Preliminary ablations showed that larger values of k introduced irrelevant context and consistently reduced faithfulness; we, therefore, adopt $k = 3$ as a conservative setting. Documents are segmented into 512-token chunks with a 50-token overlap to balance semantic coherence and retrieval granularity.

4.4 Execution Protocol and Logging

All twelve strategies are evaluated on the same set of 250 classes. Metrics are computed per sample and aggregated across the dataset. Alongside documentation quality, coverage, and faithfulness, we record end-to-end latency, token usage, and model-call counts for each strategy, enabling explicit analysis of quality–efficiency trade-offs. All prompts, hyperparameters, and evaluation settings are controlled via centralized configuration files to ensure reproducibility and consistent comparisons.

5 Results

5.1 Core Text Quality Metrics

We first examine standard text quality metrics, including lexical overlap (ROUGE-1, BLEU) and semantic similarity (BERTScore). Aggregate results are summarized in Table 1. Across architectural families, SimpleRAG (Base) achieves the highest semantic alignment, with a BERTScore of 0.598, outperforming both the Plain LLM baseline (0.595) and Self-Correction RAG (0.569). Paired Wilcoxon signed-rank tests confirm that the improvement of SimpleRAG over Plain LLM is statistically significant ($p < 0.05$). While reasoning-based prompting strategies occasionally yield marginal improvements in lexical overlap, these gains do not translate into consistent semantic improvements. Across all families, differences between Base prompting and reasoning-augmented variants (CoT, ToT, GoT) are not statistically significant with respect to BERTScore. These results suggest that increasing reasoning depth does not systematically improve semantic quality for docstring generation.

5.2 Faithfulness and Hallucination

Faithfulness, measured using an independent LLM-as-a-Judge, reveals the most pronounced differences between strategies. Quantitative scores are

| Setting | Strategy | ROUGE-1 | BLEU | BERTScore | Readability |
|-----------------|------------------|---------|-------|--------------|-------------|
| No RAG | Plain | 0.190 | 0.010 | 0.595 | 47.40 |
| | CoT | 0.191 | 0.007 | 0.602 | 38.14 |
| | GoT | 0.212 | 0.011 | 0.581 | 44.53 |
| | ToT | 0.210 | 0.012 | 0.581 | 49.56 |
| RAG | SimpleRAG | 0.200 | 0.006 | 0.598 | 39.95 |
| | CoT | 0.176 | 0.009 | 0.574 | 32.93 |
| | GoT | 0.205 | 0.008 | 0.579 | 42.86 |
| | ToT | 0.205 | 0.012 | 0.586 | 44.72 |
| RAG + Reasoning | IterCrit RAG | 0.203 | 0.009 | 0.569 | 50.68 |
| | CoT | 0.204 | 0.008 | 0.562 | 48.87 |
| | GoT | 0.221 | 0.006 | 0.585 | 50.70 |
| | ToT | 0.195 | 0.005 | 0.560 | 47.30 |

Table 1: Core text quality metrics averaged over 250 classes. Improvements in lexical and semantic similarity from reasoning strategies are inconsistent and do not correlate with faithfulness.

reported in Table 2.

Plain LLM strategies consistently cluster around a faithfulness score of 0.50, regardless of reasoning mode, indicating a persistent hallucination in the absence of grounding. In contrast, SimpleRAG (Base) achieves the highest faithfulness score (0.731), representing a 46% relative improvement over the Plain LLM baseline. This improvement is statistically significant under paired testing ($p < 0.001$).

Reasoning-augmented RAG variants do not surpass the base retrieval strategy. Although ToT-RAG and GoT-RAG achieve moderate faithfulness scores, they remain significantly below SimpleRAG ($p < 0.05$). These findings indicate that additional reasoning steps can introduce unsupported inferences rather than improve factual grounding for this task.

5.3 Documentation Coverage and Its Interaction with Faithfulness

We next analyze how comprehensively each strategy documents the code interface. Parameter, return, and exception coverage metrics are reported in Table 2, while their relationship with faithfulness is illustrated in Figure 3. SimpleRAG achieves the best balance between documentation completeness and correctness, maintaining high parameter, return, and exception coverage while also attaining the highest faithfulness. In contrast, Self-Correction RAG shows lower coverage with only moderate faithfulness, indicating that iterative critique can suppress valid interface information alongside hallucinations. Reasoning-heavy strategies occasionally increase coverage, but at a substantially higher computational cost and with-

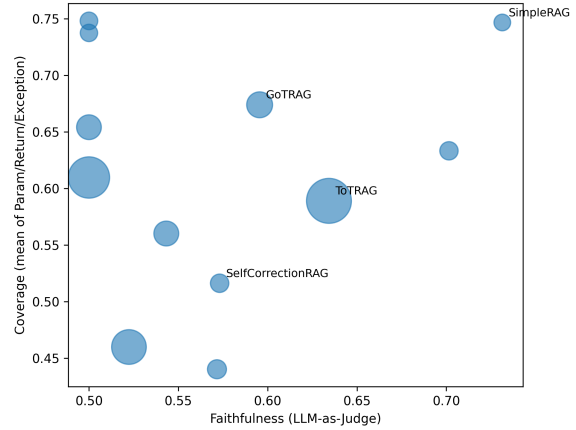


Figure 3: Relationship between documentation coverage and faithfulness. Bubble size indicates average execution time.

out corresponding gains in faithfulness. Crucially, higher coverage does not imply higher faithfulness. Encouraging exhaustive documentation can prompt models to speculate beyond the available evidence. Paired statistical tests confirm that SimpleRAG significantly outperforms Self-Correction and reasoning-enhanced variants in both coverage and faithfulness ($p < 0.05$).

5.4 Computational Cost and Efficiency

We analyze computational efficiency using execution time and its components, reported in Table 3, with the overall quality–cost trade-off shown in Figure 4. All experiments were executed on a single NVIDIA A100 GPU under identical conditions. SimpleRAG (Base) achieves a mean latency of approximately 5.9 seconds per sample, making it suitable for interactive developer workflows. In contrast, ToT-RAG incurs

| Setting | Strategy | ParamCov | ReturnCov | ExceptionCov | Faithfulness |
|-----------------|------------------|----------|--------------|--------------|--------------|
| No RAG | Plain | 0.528 | 0.925 | 0.791 | 0.501 |
| | CoT | 0.563 | 0.866 | 0.784 | 0.500 |
| | GoT | 0.463 | 0.821 | 0.679 | 0.504 |
| | ToT | 0.379 | 0.746 | 0.703 | 0.502 |
| RAG | SimpleRAG | 0.455 | 0.940 | 0.845 | 0.731 |
| | CoT | 0.489 | 0.746 | 0.664 | 0.702 |
| | GoT | 0.428 | 0.896 | 0.698 | 0.596 |
| | ToT | 0.471 | 0.657 | 0.639 | 0.634 |
| RAG + Reasoning | IterCrit RAG | 0.287 | 0.612 | 0.649 | 0.573 |
| | CoT | 0.274 | 0.508 | 0.539 | 0.572 |
| | GoT | 0.408 | 0.642 | 0.631 | 0.543 |
| | ToT | 0.267 | 0.552 | 0.560 | 0.522 |

Table 2: Documentation coverage and faithfulness scores. SimpleRAG (Base) achieves the highest faithfulness despite not maximizing coverage, indicating that reasoning-based verbosity often introduces unsupported content.

| Setting | Strategy | Latency (s) | API Calls | Retrieval (s) | Generation (s) |
|-----------------|------------------|-------------|-----------|---------------|----------------|
| No RAG | Plain | 7.44 | 1.0 | 0.00 | 7.44 |
| | CoT | 7.05 | 1.0 | 0.00 | 7.05 |
| | GoT | 21.93 | 1.0 | 0.00 | 21.93 |
| | ToT | 74.49 | 1.0 | 0.00 | 74.49 |
| RAG | SimpleRAG | 5.90 | 2.0 | 0.05 | 5.84 |
| | CoT | 8.69 | 2.0 | 0.05 | 8.64 |
| | GoT | 24.74 | 2.0 | 0.06 | 24.68 |
| | ToT | 89.52 | 2.0 | 0.28 | 89.23 |
| RAG + Reasoning | IterCrit RAG | 8.82 | 2.66 | 0.15 | 6.61 |
| | CoT | 9.89 | 2.75 | 0.20 | 6.66 |
| | GoT | 22.21 | 2.75 | 0.15 | 12.12 |
| | ToT | 50.28 | 2.69 | 0.29 | 12.05 |

Table 3: Computational cost comparison measured on a single NVIDIA A100 GPU. Reasoning-heavy strategies incur substantial latency without commensurate gains in faithfulness.

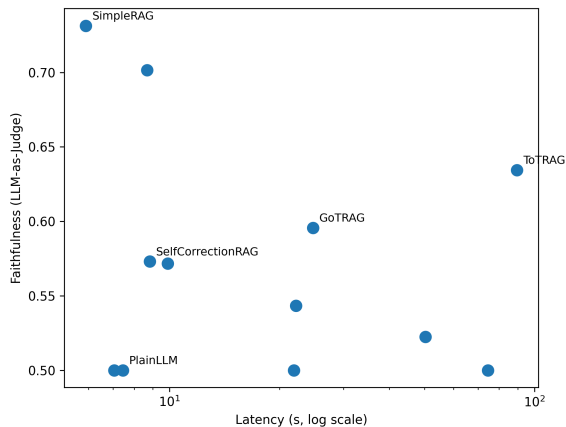


Figure 4: Trade-off between faithfulness and execution latency for all strategies. Latency is shown on a logarithmic scale.

a mean latency of approximately 89.5 seconds, representing a 15× increase over SimpleRAG. This difference is highly statistically significant ($p < 0.001$). Importantly, the increased computational cost does not yield corresponding improvements in faithfulness or semantic quality. Figure 4 shows that SimpleRAG consistently occupies the optimal region of the faithfulness–latency space, achieving the highest correctness at the lowest cost.

Overall, retrieval-based grounding consistently dominates generation-time reasoning in producing faithful docstrings. While self-correction improves over Plain LLMs, it remains unstable and inferior to SimpleRAG. Across all metrics, statistical tests confirm that SimpleRAG (Base) lies on the optimal faithfulness–cost Pareto frontier, achieving the best balance between reliability and efficiency.

| | | | |
|-----|--|---|--|
| 501 | 6 Discussion | | |
| 502 | 6.1 Why Reasoning Underperforms for Docstring Generation | | |
| 503 | | | |
| 504 | Our results show that explicit generation-time reasoning strategies—Chain-of-Thought, Tree-of-Thought, and Graph-of-Thought—do not improve faithfulness for retrieval-augmented docstring generation and often degrade it. This outcome is best explained by the nature of the task. Docstring generation is a descriptive, context-bound problem: the objective is to accurately restate information explicitly present in the code and supporting documentation. Unlike tasks that require planning or inference, documentation benefits from strict adherence to observable structure. Reasoning-based prompting expands the model’s deliberation space, encouraging speculative inferences about behavior or semantics that are not supported by the input. Without external verification, these inferences frequently lead to unsupported statements—a failure mode we term hallucination by inference. Empirically, reasoning-heavy strategies substantially increase computational cost while yielding equal or lower faithfulness than simpler alternatives. | | |
| 505 | | | |
| 506 | | | |
| 507 | | | |
| 508 | | | |
| 509 | | | |
| 510 | | | |
| 511 | | | |
| 512 | | | |
| 513 | | | |
| 514 | | | |
| 515 | | | |
| 516 | | | |
| 517 | | | |
| 518 | | | |
| 519 | | | |
| 520 | | | |
| 521 | | | |
| 522 | | | |
| 523 | | | |
| 524 | | | |
| 525 | 6.2 Retrieval as a Constraint on Generation | | |
| 526 | In contrast, retrieval-based grounding acts as an effective constraint on generation. By injecting concrete external context, SimpleRAG narrows the output space to statements supported by retrieved evidence, reducing hallucination without additional control mechanisms. Even lightweight retrieval configurations achieve large gains in faithfulness, indicating that retrieval actively regularizes generation rather than merely supplementing parametric knowledge. | | |
| 527 | | | |
| 528 | | | |
| 529 | | | |
| 530 | | | |
| 531 | | | |
| 532 | | | |
| 533 | | | |
| 534 | | | |
| 535 | | | |
| 536 | 6.3 Implications for Agentic and Self-Correction Frameworks | | |
| 537 | | | |
| 538 | The performance of Self-Correction RAG highlights important limitations of agentic approaches for documentation tasks. Although iterative critique improves upon Plain LLM baselines, it remains consistently inferior to SimpleRAG and exhibits high variance. Because the critic and generator share the same incomplete context and inductive biases, self-correction often reinforces earlier assumptions rather than correcting them, aligning with prior findings that effective agentic correction requires external verification signals. Overall, our results indicate that for structured, context-bound | | |
| 539 | | | |
| 540 | | | |
| 541 | | | |
| 542 | | | |
| 543 | | | |
| 544 | | | |
| 545 | | | |
| 546 | | | |
| 547 | | | |
| 548 | | | |
| 549 | | | |
| | | tasks such as code documentation, grounding mechanisms are more reliable than internal deliberation or self-reflection. | 550 551 552 |
| | 7 Conclusion | | 553 |
| | This work highlights the importance of task-aware architectural design in large language model-based software engineering systems, showing that increasingly complex reasoning mechanisms do not universally improve reliability. We present a systematic empirical evaluation of retrieval, reasoning, and self-correction strategies for automated Python docstring generation across twelve configurations and multiple quality, coverage, faithfulness, and cost metrics. | | 554 555 556 557 558 559 560 561 562 563 |
| | Our results demonstrate that retrieval-based grounding is the primary driver of reliable documentation, while generation-time reasoning yields diminishing—and, in several cases, negative—returns. A lightweight Retrieval-Augmented Generation (RAG) baseline consistently outperforms both standalone LLMs and reasoning-heavy pipelines in faithfulness while incurring substantially lower computational costs. In contrast, explicit reasoning strategies—including Chain-of-Thought, Tree-of-Thought, and Graph-of-Thought—significantly increase inference latency without improving factual correctness and often introduce additional unsupported content. We characterize this failure mode as hallucination by inference, where reasoning encourages speculative abstraction that is misaligned with the extractive, context-bound nature of documentation tasks. | | 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 |
| | These findings have direct implications for the design of automated documentation systems. Rather than indiscriminately applying complex agentic pipelines, retrieval-first architectures offer a more effective, efficient, and sustainable solution for developer-facing tools, where latency and reliability are critical constraints. Future work will extend this analysis to additional programming languages, documentation formats, and codebase-aware retrieval settings, and explore adaptive systems that invoke reasoning only when ambiguity is detected. We hope this study encourages more principled evaluation of when reasoning is necessary—and when efficient grounding alone is sufficient—in practical software engineering applications. | | 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 |

598 Limitations

599 These results should not be interpreted as evi-
600 dence against the usefulness of reasoning strate-
601 gies in tasks that genuinely require planning, ex-
602 ecution, or multi-step inference. Our findings
603 are specific to Python class-level docstring gen-
604 eration, a static and context-bound documentation
605 task where correctness depends primarily on faith-
606 ful restatement rather than inference. Several ad-
607 ditional limitations remain. First, our evaluation
608 focuses on a single programming language and
609 documentation style, and results may not directly
610 generalize to other languages or documentation
611 formats. Second, although the dataset spans di-
612 verse GitHub repositories, its size is modest, and
613 larger-scale studies could further validate the ob-
614 served trends. Third, faithfulness is assessed us-
615 ing an LLM-as-a-Judge, which—while effective
616 for relative comparison—may reflect biases of the
617 evaluator model. Finally, we do not include human
618 evaluation, which could provide complementary
619 insights into developer-perceived usefulness and
620 clarity.

621 Ethics Statement

622 We acknowledge two primary ethical considera-
623 tions in this work. First, AI-assisted tools were
624 used during writing and code development solely to
625 improve clarity and presentation. All experimental
626 design, analysis, and conclusions were developed
627 and verified by the human authors, who take full
628 responsibility for the content.

629 Second, all experiments were conducted using
630 publicly available, open-source code and datasets
631 under their respective licenses. No private, sen-
632 sitive, or personally identifiable data were used,
633 and data ownership and licensing constraints were
634 respected.

635 More broadly, our findings support the re-
636 sponsible and sustainable use of large language
637 models. By showing that lightweight retrieval-
638 based grounding achieves higher reliability than
639 reasoning-heavy pipelines while reducing infer-
640 ence costs by up to 15×, this work promotes ef-
641 ficient deployment and aligns with the principles
642 of Green AI and responsible model usage.

643 References

644 Anna Bavaresco, Raffaella Bernardi, Leonardo Berto-
645 lazzi, Desmond Elliott, Raquel Fernández, Albert

Gatt, Esam Ghaleb, Mario Giulianelli, Michael
646 Hanna, Alexander Koller, Andre Martins, Philipp
647 Mondorf, Vera Neplenbroek, Sandro Pezzelle, Bar-
648 bara Plank, David Schlangen, Alessandro Suglia,
649 Aditya K Surikuchi, Ece Takmaz, and Alberto
650 Testoni. 2025. [LLMs instead of human judges? a
651 large scale empirical study across 20 NLP evalua-
652 tion tasks](#). In *Proceedings of the 63rd Annual Meet-
653 ing of the Association for Computational Linguistics
654 (Volume 2: Short Papers)*, pages 238–255, Vienna,
655 Austria. Association for Computational Linguistics. 656

Maciej Besta, Nils Blach, Ales Kubicek, Robert Ger-
657 stenberger, Michał Podstawski, Lukas Gianinazzi,
658 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadom-
659 ski, Piotr Nyczyk, and Torsten Hoeffler. 2024. [Graph
660 of thoughts: solving elaborate problems with large
661 language models](#). In *Proceedings of the Thirty-
662 Eighth AAAI Conference on Artificial Intelligence
663 and Thirty-Sixth Conference on Innovative Applica-
664 tions of Artificial Intelligence and Fourteenth Sym-
665 posium on Educational Advances in Artificial Intelli-
666 gence, AAAI’24/IAAI’24/EAAI’24*. AAAI Press. 667

Nicola Dainese, Alexander Ilin, and Pekka Marttinen.
668 2024. [Can docstring reformulation with an llm im-
669 prove code generation?](#) In *Conference of the Euro-
670 pean Chapter of the Association for Computational
671 Linguistics*. 672

Sherif Elmitwalli, John Mehegan, Sophie Braznell, and
673 Allen Gallagher. 2025. [Scalable evaluation frame-
674 work for retrieval augmented generation in tobacco
675 research using large language models](#). *Scientific Re-
676 ports*, 15. 677

Max Grusky. 2023. [Rogue scores](#). In *Proceedings
678 of the 61st Annual Meeting of the Association for
679 Computational Linguistics (Volume 1: Long Papers)*,
680 pages 1914–1934, Toronto, Canada. Association for
681 Computational Linguistics. 682

Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas
683 Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-
684 Yu, Armand Joulin, Sebastian Riedel, and Edouard
685 Grave. 2022. [Atlas: Few-shot learning with retrieval
686 augmented language models](#). *arXiv preprint arXiv:
687 2208.03299*. 688

Patrick Lewis, Ethan Perez, Aleksandara Piktus,
689 F. Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
690 rich Kuttler, M. Lewis, Wen tau Yih, Tim Rock-
691 täschel, Sebastian Riedel, and Douwe Kiela. 2020.
692 Retrieval-augmented generation for knowledge-
693 intensive nlp tasks. *Neural Information Processing
694 Systems*. 695

Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad
696 Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-
697 tacherjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu,
698 Kai Shu, Lu Cheng, and Huan Liu. 2025a. [From gen-
699 eration to judgment: Opportunities and challenges of
700 LLM-as-a-judge](#). In *Proceedings of the 2025 Con-
701 ference on Empirical Methods in Natural Language
702 Processing*, pages 2757–2791, Suzhou, China. Asso-
703 ciation for Computational Linguistics. 704

| | | | |
|-----|---|--|-----|
| 705 | Feiyang Li, Peng Fang, Zhan Shi, Arijit Khan, | Yitong Liu, Si Zheng, and Xiangke Liao. 2025b. Un- | 762 |
| 706 | Fang Wang, Weihao Wang, Zhangxin-hw, and Cui | seen horizons: Unveiling the real capability of llm | 763 |
| 707 | Yongjian. 2025b. CoT-RAG: Integrating chain of | code generation beyond the familiar. In <i>Proceed-</i> | 764 |
| 708 | thought and retrieval-augmented generation to en- | ings of the <i>IEEE/ACM 47th International Conference</i> | 765 |
| 709 | hance reasoning in large language models. In <i>Find-</i> | on <i>Software Engineering, ICSE '25</i> , page 604–615. | 766 |
| 710 | <i>ings of the Association for Computational Linguistics:</i> | IEEE Press. | 767 |
| 711 | <i>EMNLP 2025</i> , pages 3119–3171, Suzhou, China. As- | | |
| 712 | sociation for Computational Linguistics. | | |
| 713 | Shangqing Liu, Yu Chen, Xiaofei Xie, Jing Kai Siow, | Yuxin Zhang, Yuxia Zhang, Zeyu Sun, Yanjie Jiang, and | 768 |
| 714 | and Yang Liu. 2021. Retrieval-augmented generation | Hui Liu. 2025c. Laura: Enhancing code review gen- | 769 |
| 715 | for code summarization via hybrid GNN . In <i>Internat-</i> | eration with context-enriched retrieval-augmented | 770 |
| 716 | <i>ional Conference on Learning Representations</i> . | llm. <i>Preprint</i> , arXiv:2512.01356. | 771 |
| 717 | Kishore Papineni, Salim Roukos, Todd Ward, and Wei- | | |
| 718 | Jing Zhu. 2002. Bleu: a method for automatic evalu- | | |
| 719 | ation of machine translation. In <i>Proceedings of the</i> | | |
| 720 | <i>40th Annual Meeting on Association for Computa-</i> | | |
| 721 | <i>tional Linguistics, ACL '02</i> , page 311–318, USA. | | |
| 722 | Association for Computational Linguistics. | | |
| 723 | Gireesh Sundaram, Balaji Venkatesh V, and Sundharaku- | | |
| 724 | mar K B. 2025. Docstringeval: Evaluating the ef- | | |
| 725 | fectiveness of language models for code explanation | | |
| 726 | through docstring generation . In <i>2025 International</i> | | |
| 727 | <i>Conference on Emerging Technologies in Computing</i> | | |
| 728 | <i>and Communication (ETCC)</i> , pages 1–7. | | |
| 729 | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten | | |
| 730 | Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, | | |
| 731 | and Denny Zhou. 2022. Chain-of-thought prompt- | | |
| 732 | ing elicits reasoning in large language models . In | | |
| 733 | <i>Proceedings of the 36th International Conference on</i> | | |
| 734 | <i>Neural Information Processing Systems, NIPS '22</i> , | | |
| 735 | Red Hook, NY, USA. Curran Associates Inc. | | |
| 736 | Guang Yang, Yu Zhou, Wei Cheng, Xiangyu Zhang, | | |
| 737 | Xiang Chen, Terry Yue Zhuo, Ke Liu, Xin Zhou, | | |
| 738 | David Lo, and Taolue Chen. 2025. Less is more: | | |
| 739 | Docstring compression in code generation . <i>ACM</i> | | |
| 740 | <i>Trans. Softw. Eng. Methodol.</i> Just Accepted. | | |
| 741 | Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, | | |
| 742 | Thomas L. Griffiths, Yuan Cao, and Karthik | | |
| 743 | Narasimhan. 2023. Tree of thoughts: deliberate prob- | | |
| 744 | lem solving with large language models . In <i>Pro-</i> | | |
| 745 | <i>ceedings of the 37th International Conference on</i> | | |
| 746 | <i>Neural Information Processing Systems, NIPS '23</i> , | | |
| 747 | Red Hook, NY, USA. Curran Associates Inc. | | |
| 748 | Haozhen Zhang, Tao Feng, and Jiaxuan You. 2025a. | | |
| 749 | Graph of records: Boosting retrieval augmented gen- | | |
| 750 | eration for long-context summarization with graphs . | | |
| 751 | In <i>Proceedings of the 63rd Annual Meeting of the</i> | | |
| 752 | <i>Association for Computational Linguistics (Volume 1:</i> | | |
| 753 | <i>Long Papers)</i> , pages 23780–23799, Vienna, Austria. | | |
| 754 | Association for Computational Linguistics. | | |
| 755 | Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. | | |
| 756 | Weinberger, and Yoav Artzi. 2020. Bertscore: Eval- | | |
| 757 | uating text generation with bert . In <i>International</i> | | |
| 758 | <i>Conference on Learning Representations</i> . | | |
| 759 | Yuanliang Zhang, Yifan Xie, Shanshan Li, Ke Liu, | | |
| 760 | Chong Wang, Zhouyang Jia, Xiangbing Huang, Jie | | |
| 761 | Song, Chaopeng Luo, Zhizheng Zheng, Rulin Xu, | | |

A.1 Output of generated docstrings:

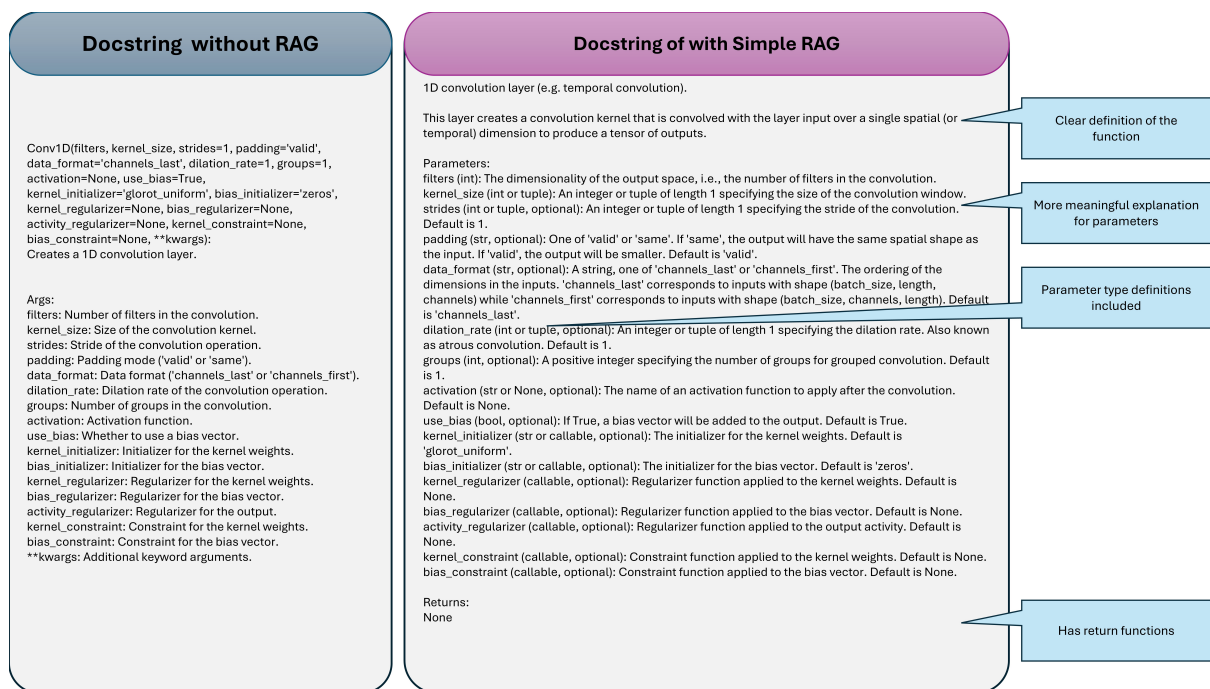


Figure 5: This image shows two docstrings generated. The one on the left shows the docstring generated without RAG and one on the right shows with RAG.

A.2 System Prompts:

You are an expert Python programmer and documentation specialist with deep knowledge of PEP 257 standards and Google-style docstring formatting.

****Your Role:****

- Generate comprehensive, accurate docstrings following PEP 257 standards
- Use Google-style docstring format with proper sections
- Include type hints and detailed parameter descriptions
- Document exceptions and return values thoroughly

****Context Usage:****

- Use retrieved context only when directly relevant to the code
- Incorporate context insights to enhance accuracy
- Avoid including irrelevant information from context

****Output Requirements:****

- Return **ONLY** the docstring content starting with triple quotes
- Do **NOT** include the original Python code
- Follow proper indentation and formatting

Figure 6: System Prompt

Generate a comprehensive Python docstring following PEP 257 standards and Google-style formatting.

****CRITICAL REQUIREMENT:****

- Return ONLY the docstring content (the text between triple quotes)
- Do NOT include the original Python code
- Do NOT include the function/class definition

****Context Integration:****

- Use any relevant context provided earlier to enhance accuracy
- Incorporate context insights into parameter descriptions

****Code to document:****

```
```python  
{code}
```

Figure 7: Final Generation Prompt

Generate a comprehensive search query to find relevant documentation for creating Python docstrings.

**\*\*Code Analysis:\*\***

Analyze the provided Python code to extract:

- Function/class name and purpose
- Parameter names and types
- Return type and value

**\*\*Search Strategy:\*\***

Create a query that will retrieve:

- PEP 257 documentation and examples
- Similar function documentation patterns

**\*\*Code to analyze:\*\***

{code}

**\*\*Generate a focused search query for this code.\*\***

Figure 8: Retrieval Context Query Generation Prompt

```
Generate a Python docstring for the code below.

Instructions:
1. **Think Step-by-Step**: First, analyze the code logic, parameters, returns, and exceptions. Explain your reasoning clearly.
2. **Generate Docstring**: Based on your analysis, generate the final docstring.

Format:
[REASONING]
<Your step-by-step analysis here>
[/REASONING]

[DOCSTRING]
<Your final docstring here (triple-quoted)>
[/DOCSTRING]

Code:
{code}
```

Figure 9: Chain of Thought Prompt

Decompose the task of generating a docstring for the following Python class into multiple alternative plans.

- Generate three distinct docstring drafts focusing on different aspects of the class.
- Evaluate which draft best covers purpose, parameters, returns, and exceptions.
- Synthesize the best elements into a final PEP-style docstring.

Avoid unsupported assumptions or hallucinated details.

```
Code:
{code}
```

Figure 10: Tree of Thought Prompt

Independently analyze the following aspects of the Python class:

- Class purpose
- Parameters and attributes
- Return behavior
- Exceptions

Combine these analyses into a single coherent PEP-style docstring. Ensure all statements are directly supported by the code or retrieved context.

```
Code:
{code}
```

Figure 11: Graph of Thought Prompt

Evaluate the quality of the generated Python docstring against the original code using comprehensive criteria.

**\*\*Evaluation Framework:\*\***

1. **\*\*Completeness\*\***: Documents all parameters/returns/exceptions
2. **\*\*Accuracy\*\***: Descriptions match implementation
3. **\*\*Formatting\*\***: PEP 257 compliance
4. **\*\*Clarity\*\***: Professional language

**\*\*Scoring:\*\***

- EXCELLENT (7-8 points)
- GOOD (5-6 points)
- NEEDS\_IMPROVEMENT (0-4 points)

**\*\*Code to evaluate:\*\***

{code}

**\*\*Generated docstring:\*\***

{initial\_docstring}

**\*\*Assessment (EXCELLENT/GOOD/NEEDS\_IMPROVEMENT):\*\***

Figure 12: Iterative Critique RAG – Critique Prompt

```

You are a strict technical judge evaluating Python
docstrings.

Context from Knowledge Base:
{retrieved_context}

Generated Docstring:
{generated_docstring}

Evaluation Rubric:
1. Hallucination Check: Does the docstring claim
parameters/returns not present in the code or context?
2. Contradiction Check: Does it contradict the provided
context logic?
3. Support Check: Is the description supported by the
context?

Task:
Assign a Faithfulness Score from 0.0 to 1.0.
- 1.0: Fully supported by context/code, no hallucinations.
- 0.5: Partially supported, some generic descriptions.
- 0.0: Major hallucinations or contradictions.

Return ONLY the numeric score in this format: Score:
<number>

```

784

Figure 13: LLM-as-a-Judge Prompt

#### A.4 Evaluation Framework:

785

To ensure a holistic assessment of the generated docstrings, we employed a comprehensive suite of evaluation metrics categorized into two main areas: code evaluation metrics and documentation coverage metrics.

786

787

788

##### A.4.1 Code evaluation metrics

789

These metrics assess the fundamental quality of the generated natural language text. We used 5 key metrics: ROUGE-1 score, BLEU score, BERT score, ease of reading, and conciseness. For measuring accuracy, the main metric that we use is the BERT score. The BERT score is calculated as follows:

790

791

792

$$BERT = 2 \times \frac{P_{BERT} \times R_{BERT}}{P_{BERT} + R_{BERT}}$$

Where  $P_{BERT}$  is the Precision BERT score, and  $R_{BERT}$  is the Recall BERT score.

793

794

To calculate the ease of reading, we use the Flesch reading ease metric, which is calculated based on the ratio of total words to total syllables, as mentioned below:

$$Clarity = \frac{206.825 - 1.5 \frac{Total_w}{Total_s} - 84.6 \frac{Total_{sy}}{Total_w}}{100}$$

where  $w$  are words,  $s$  are sentences, and  $sy$  is the generated syllabus.

Conciseness measures how precise the docstring is.

$$Conciseness = \frac{Compress_{Generated\ text}}{Compress_{Ground\ truth}}$$

The ideal conciseness score will be around 1, which indicates that the compressed text of the ground truth and the generated text are of the same length.

#### A.4.2 Documentation Coverage Metrics

These metrics evaluate how well the docstring documents the functional interface of the code.

Parameter Coverage measures the proportion of function and method parameters that are mentioned in the generated docstring vs. the actual code.

$$Parameter\ coverage = \frac{D_{Parameters}}{C_{Parameters}}$$

where  $D_{Parameters}$  are the generated docstring parameters, and  $C_{Parameters}$  are the code docstring parameters.

Return coverage is a binary score that checks whether a return statement in the code is documented.

$$Return\ coverage = \begin{cases} 1 & \text{if code has docstring and no return statement} \\ 1 & \text{if code and docstring both have a return statement} \\ 0 & \text{if code has a return statement but the docstring does} \\ & \text{not, or vice versa} \end{cases} \quad (1)$$

Faithfulness scores represent the relationship between the tokens in the generated docstrings and the retrieved context.

$$Faithfulness\ score = \frac{D_{Tokens} \cap RC_{Tokens}}{D_{Tokens}}$$

where  $D_{Tokens}$  represents tokens generated by the docstring, and  $RC_{Tokens}$  represents the retrieved context.

Exception Score measures the proportion of explicitly raised exceptions in the code that are mentioned in the docstring.

$$Exception\ coverage = \frac{D_{Exception}}{C_{Exceptions}}$$

where  $D_{Exceptions}$  stands for the number of exceptions raised in the generated docstring, and  $C_{Exceptions}$  is the number of exceptions in the code.