

Code for Machines, Not Just Humans: Quantifying AI-Friendliness with Code Health Metrics

Markus Borg
CodeScene and Lund University
Malmö, Sweden
markus.borg@codescene.com

Adam Tornhill
CodeScene
Malmö, Sweden
adam.tornhill@codescene.com

Nadim Hagatulah
Lund University
Lund, Sweden
nadim.hagatulah@cs.lth.se

Emma Söderberg
Lund University
Lund, Sweden
emma.soderberg@cs.lth.se

Abstract

We are entering a hybrid era in which human developers and AI coding agents work in the same codebases. While industry practice has long optimized code for human comprehension, it is increasingly important to ensure that LLMs with different capabilities can edit code reliably. In this study, we investigate the concept of “AI-friendly code” via LLM-based refactoring on a dataset of 5,000 Python files from competitive programming. We find a meaningful association between CodeHealth, a quality metric calibrated for human comprehension, and semantic preservation after AI refactoring. Our findings confirm that human-friendly code is also more compatible with AI tooling. These results suggest that organizations can use CodeHealth to guide where AI interventions are lower risk and where additional human oversight is warranted. Investing in maintainability not only helps humans; it also prepares for large-scale AI adoption.

CCS Concepts

• **Software and its engineering** → **Software maintenance tools.**

Keywords

software maintainability, code quality, refactoring, AI assistants

ACM Reference Format:

Markus Borg, Nadim Hagatulah, Adam Tornhill, and Emma Söderberg. 2026. Code for Machines, Not Just Humans: Quantifying AI-Friendliness with Code Health Metrics. In . ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

For decades, the maxim has been that “*programs must be written for people to read, and only incidentally for machines to execute*” [2]. Human-readable code is essential for maintaining secure, reliable, and efficient software development [11, 30, 37]. But with the advent

of AI-assisted coding, source code now has a broader audience: machines need to understand its intent, too.

Early field observations collected by Thoughtworks¹ suggest that AI-assisted coding tools perform better on high-quality code. In particular, well-factored modular code seems to reduce the risk of hallucination and lead to more accurate suggestions. They refer to this as “*AI-friendly code design*” in their April 2025 Technology Radar, and discuss how established best practices so far align with AI-friendliness [43]. If this observation holds, it implies that code optimized for human comprehension is also easier for Large Language Models (LLMs) to process and evolve.

The implications of AI-friendliness are profound: as of 2025, about 80% of developers already use AI tools in their work [31], with adoption projected to grow rapidly. Gartner predicts an increase from 14% in early 2024 to 90% by 2028 [46]. At the same time, less than 10% of organizations have been reported to methodically track technical debt [27], and developers have been found to waste up to 42% of their time due to poor code quality [42]. These figures suggest that much of today’s production code may be structurally unfit for reliable AI intervention, increasing the risk of bugs and expensive rework.

In this paper, we investigate the relationship between code quality and AI-friendliness. We measure quality using the *CodeHealth* (CH) metric, which has been validated as predictive of defects and development effort in previous studies [6, 9, 44]. The metric has also been used in recent industry-facing AI studies [8, 29]. We use the success rate of AI-generated refactoring, i.e., improving the design of existing code without changing its behavior, as a proxy for AI-friendliness. Refactoring lets us use passing unit tests as an oracle for functional correctness. Thus, an AI refactoring is *correct* if tests pass and *beneficial* if CH increases.

Our results confirm that human-friendly code is more compatible with AI tooling. LLMs tasked with refactoring have significantly lower break rates on code in the Healthy CodeHealth range ($CH \geq 9$), with corresponding risk reductions of 15-30%. Furthermore, we show that CH outperforms *perplexity* (PPL, an LLM-intrinsic confidence metric) and Source Lines of Code (SLOC) as a predictor of refactoring correctness.

These findings can support software organizations in the adoption of AI-assisted coding. We propose using CH to identify parts of the code that are ready for AI processing, as well as highlighting

¹<https://www.thoughtworks.com/>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference’17, Washington, DC, USA

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

code with too much risk of breaking. More broadly, this research adds a missing piece to AI adoption: a shared code-quality metric that aligns humans and machines. While earlier research positioned code quality as a business imperative, we posit that code quality is a prerequisite for safe and effective use of AI – which might prove existential for software organizations in the next decade.

The remainder of this paper is organized as follows. Section 2 describes background and related work. In Section 3, we introduce our research questions, dataset, and method. Section 4 reports our results. Finally, Section 5 discusses the implications before Section 6 concludes with directions for future work.

2 Background and Related Work

This section first introduces CH and PPL, then presents related work on code comprehension and AI refactoring.

2.1 Maintainability and CodeHealth

CodeHealth™ (CH) is a quality metric used in the CodeScene software engineering intelligence platform. Its goal is to capture how cognitively difficult it is for human developers to comprehend code. CodeScene identifies *code smells* [25], e.g., God classes, deeply nested logic, and duplicated code. For Python, which we target in this paper, CodeScene detects 25 code smells.

CodeScene combines the number and severity of detected smells into a file-level score from 1 to 10. Lower scores indicate higher cognitive load for humans, i.e., higher maintenance effort. CodeScene categorizes files as belonging to one of three CH intervals: *Healthy* ($CH \geq 9$), *Warning* ($4 \leq CH < 9$), and *Alert* ($CH < 4$). In this study, we refer to both Warning and Alert files as *Unhealthy*.

We have previously validated the CH metric in a series of studies. In a study on the manually annotated *Maintainability Dataset* [38], we reported that CH aligns better with human maintainability judgments than competing metrics and the average human expert [6]. Furthermore, we have validated the metric from a business perspective through the association between CH on the one hand and file-level defect density and development time on the other hand [9, 44]. Building on our previous work, we now investigate the relation between CH and AI-friendliness.

2.2 LLM Confidence and Perplexity

Several metrics have been proposed to assess the confidence of LLM output. Internal metrics such as entropy, mutual information, and PPL can potentially be used as proxies of output correctness and hallucination risk. In the software domain, Sharma *et al.* showed that higher entropy and mutual information correlate with lower functional correctness of the generated code [39]. Earlier work by Ray *et al.* also showed that code containing defects had higher entropy than correct code [34], which in turn builds on the foundational work on the naturalness of code by Hindle *et al.* [19].

Ppl, originally proposed for speech recognition tasks [21], is the exponential of the prediction model’s average entropy over a sequence. Considering entropy to represent “average surprise,” PPL rather turns this into a scale of “how many choices.” It is used as a measure of how confident an LLM is in predicting the next token based on the previous tokens – a higher score means lower confidence. Mathematically, let $t_{1:N}$ be a token sequence. The

average cross-entropy (H) is:

$$H = -\frac{1}{N} \sum_{i=1}^N \log p(t_i | t_{<i}).$$

PPL is $\exp(H)$. In this work, we get PPL directly from the five Hugging Face models under study [13].

2.3 Code Comprehension and Perplexity

Gopstein *et al.* coined the term *atoms of confusion* to refer to small, isolated patterns in code that confuse humans [14]. An atom is the smallest unit of code that can cause confusion. Through controlled experiments, they identified a set of atoms that significantly increased the rate of misunderstandings. Examples include “Assignment as Value,” e.g., `V1 = V2 = 3`; and “Logic as Control Flow,” e.g., `V1 && F2()`;

Building on the atoms concept, Abdelsalam *et al.* discussed the risks of comprehension differences between human programmers and LLMs in the current hybrid development landscape [1]. They argue that if humans and LLMs are confused by different characteristics in code, there is a risk of misalignment. By comparing LLM PPL and EEG responses, they found that LLMs and humans struggle with similar issues in code. A previous controlled experiment by Casulnuovo *et al.*, also suggests a connection between human comprehension and how surprising a code snippet is to LLMs [12].

Kotti *et al.* explored the PPL of different programming languages [23]. Based on 1,008 files from GitHub projects representing different languages, they found systematic differences between strongly-typed and dynamically typed languages. For example, across LLMs, they found that Perl resulted in high PPL and Java in low. They speculate that their PPL findings could be used to assess the suitability of LLM-based code completion in specific projects. In this study, we explore the PPL concept for Healthy vs. Unhealthy Python code.

2.4 AI Refactoring

Research in automated code refactoring has, like so many software engineering applications, been disrupted by LLMs. Several studies investigate the potential of various combinations of LLMs and prompts to improve design and remove code smells effectively. Pomian *et al.* introduced the tool EM-Assist [33] to automatically suggest and perform extract method refactorings. The same team has continued working on the more sophisticated refactoring operation *move method* [4]. Tornhill *et al.* have developed the tool ACE to automatically remove five CodeScene code smells [45].

During 2025, agentic AI has been a major trend in industry and research. Coding agents are driven by LLMs, but they typically claim to 1) understand codebases beyond limited context windows and to 2) maintain a memory. This enables agents to take on larger tasks and operate more autonomously. The arguably most popular agent in industry at this time of writing is Anthropic’s Claude Code, carefully described by Watanabe *et al.* [48]. We refer to reviews by He *et al.* [17] and Wang *et al.* [47] for contemporary overviews of the academic literature.

Refactoring is one of the many tasks investigated for coding agents. For example, Xu *et al.* presented MANTRA, a multi-agent framework for refactoring [49]. MANTRA organizes three agents into 1) *developer*, 2) *reviewer*, and 3) *repair* roles for the refactoring task, and outperforms direct LLM-usage. However, Claude is also a

capable refactoring agent despite not being a multi-agent solution. In this study, we study Claude (v2.0.13) as a representative example of contemporary agentic AI refactoring. Moreover, MANTRA is not publicly available at the time of this writing.

3 Method

Our goal is to explore how code characteristics influence the capabilities of AI refactoring, with a particular focus on CH. We formulate three Research Questions (RQs) that explore code characteristics in light of AI-friendliness. First, we take an LLM-intrinsic perspective and study PPL. This connects to related work and provides a baseline. Second, we look at how the CH is associated with the success rate of a downstream AI task. Third, we examine the predictive power of CH compared to SLOC and PPL.

- RQ₁** How does perplexity differ between Healthy and Unhealthy code?
- RQ₂** How does the AI refactoring break rate differ between Healthy and Unhealthy code?
- RQ₃** To what extent can CodeHealth predict the AI refactoring break rate?

3.1 Dataset Creation

We sample from the CodeContests dataset hosted on GitHub² by Google DeepMind. The dataset was introduced by Li *et al.* [26] as training data for AlphaCode and contains more than 12 million solutions (correct and incorrect) to competitive programming problems from five sources. We study code in this domain because the problems come with carefully crafted test cases that verify functional correctness, providing a practical oracle after refactoring.

We construct a dataset of 5,000 solutions based on four design choices. First, we decided to focus on solutions written in Python to increase novelty and convenience. While CodeContests also contains solutions in Java and C++, Java has been extensively studied in refactoring research, and C++ has a more complex compile-and-test procedure. Second, we require at least one CodeScene code smell in the solutions. Removing code smells is a realistic refactoring goal. Third, we chose to only study solutions containing between 60 and 120 SLoC. There is a strong correlation between maintainability and size [41], thus we control for this, at least partly, already during the sampling. Fourth, we actively seek diversity in the dataset, as many solutions are highly similar. We use CodeBleu [35] for similarity calculations using the following weights: `n-gram=0.1`, `weighted-n-gram=0.4`, `ast-match=0.5`, and `dataflow-match=0`.

The practical sampling process followed these steps (number of remaining solutions in parentheses):

- (1) Download the full CodeContests dataset (13,210,440).
- (2) Filter to only Python 3 solutions (1,502,532).
- (3) Remove all identical (Type 1 clone) solutions (1,390,531).
- (4) Remove all solutions with no CodeScene code smells (88,156).
- (5) Strip comments and unbound string literals.
- (6) Filter to solutions with 60–120 SLOC (18,074).
- (7) Partition into two strata: Healthy and Unhealthy code³.
- (8) For each stratum, repeat until 2,500 solutions are included:

- (b) Run its corresponding test cases; skip if any test case fails.
- (c) Compute CodeBleu similarity to existing samples; skip if similarity ≥ 0.9 .
- (d) Add the solution to the stratum.
- (9) Merge the two strata to constitute the final dataset (5,000).

3.2 Selection of Large Language Models

We select six LLMs for evaluation. Five are open-weight models with about 20-30B parameters, runnable on our local datacenter. We select four models based on popularity and download statistics from Hugging Face and complement them with an LLM recently published by IBM – we refer to these as *medium-sized LLMs*. Moreover, we include a State-of-the-Art (SotA) LLM that we prompt using Anthropic’s API.

Gemma gemma-3-27b-it (Google) – Mar 2025.

GLM GLM-4-32B-0414 (Zhipu AI) – Apr 2025.

Granite Granite-4.0-H-Small (IBM) – Oct 2025.

GPT gpt-oss-20b (OpenAI) – Aug 2025.

Qwen Qwen3-Coder-30B-A3B-Instruct (Alibaba Cloud) – Aug 2025.

Sonnet claude-sonnet-4-5-20250929 (Anthropic) – Sep 2025.

We set the sampling temperature to 0.7 for all LLMs to enable refactoring diversity under a uniform setting and we cap generation to a maximum of 8,192 new tokens. All other settings use defaults.

3.3 RQ₁ Perplexity

We extract PPL scores for all 5,000 samples for the five LLMs under study. Then we split the samples into Healthy and Unhealthy and state the following null hypothesis for each LLM:

H_0 : The perplexity distributions for Healthy and Unhealthy code are identical.

Descriptive statistics revealed a small set of outliers with extremely high PPL. Manual analysis showed specific patterns that substantially inflate the LLMs’ next-token uncertainty. In competitive programming, some contestants hard-code astronomically big integers or very long strings for later use in the solutions. To mitigate this, we applied one-sided robust z-score filtering and removed samples with $z > 2.5$ from the output from the LLMs (about 5%).

PPL scores are typically right-skewed, and we assess normality using the Shapiro–Wilk test. Even after filtering upper-tail outliers, we rejected normality for all LLM outputs. Subsequently, we rely on non-parametric Mann-Whitney U tests for significance testing. Finally, we use Holm correction (Holm-Bonferroni step-down) to adjust for multiple comparisons across the five models (family-wise $\alpha = 0.05$). We report two-sided Holm-corrected p-values and Cliff’s δ as the effect size.

3.4 RQ₂ Refactoring Break Rate

We tasked all medium-sized LLMs to refactor the 5,000 samples. For Sonnet, which has a high token price, we refactor 1,000 random samples equally split between Healthy and Unhealthy. For each combination of LLM and sample, we used the same prompt following a general structure, designed to be generic:

- (A) Context in the form of a role description to steer the model toward the right parts of its knowledge.
- (B) A concrete task for the model to perform.

²https://github.com/google-deeppmind/code_contests

³A minor threshold-specification issue in the sampling script treated CH = 9 as Unhealthy. All analyses in this paper use the intended convention (Healthy: ≥ 9).

(a) Sample a random candidate solution.

- (C) Instructions for how the model should format the response.
 (D) Input data for which the task shall be completed.

Refactoring prompt

^(A)Act as an expert software engineer. ^(B)Your task is to refactor the following Python code for maintainability and clean code. ^(C) Respond ONLY with the complete, refactored Python code block. Do not add any explanations, comments, or introductory sentences.
^(D)Original code to refactor: ``` python <CODE> ```

We complemented the LLMs with Claude (v2.0.13) to also investigate a SotA agentic approach to refactoring. We selected Claude specifically because it is considered industry-leading and currently tops public SWE-bench results [22]. While Claude operates in an interactive mode in the terminal, we controlled the environment and presented the task in a manner consistent with the LLM setup. Claude is a costly service, and we decided to target the same 1,000 samples as for Sonnet. We aimed at a final cost of less than \$100 for Anthropic's solutions.

We organized batches of random samples in separate folders to mimic how a developer might work, starting with a small sample of 20 files for a pilot run – followed by 200, 400, and 380. We provided the instructions (A–C) in a CLAUDE.md file in the common parent folder, appended with: "Each file in this folder is independent; you do not have to worry about dependencies between them."

For each sample folder used with Claude, we added a configuration file (.claude/settings.json) to pin the model version and restrict the agent's tool use. Specifically, we set the model to claude-sonnet-4-5-20250929, explicitly disabled Bash, WebFetch, and WebSearch, and also disabled any MCP use.

For each refactoring session, we followed these steps:

- (1) Launch Claude in the current folder.
- (2) Answer yes to "Do you trust the files in this folder?"
- (3) Ask Claude "Can you see CLAUDE.md in the parent folder?" and grant read permission.
- (4) Give the instruction: "Refactor the Python files in this folder for maintainability."
- (5) When Claude requests edit permission, approve editing of all files for the session.

After each single refactoring pass, we recorded descriptive statistics of the output code, calculated its CH, and executed the corresponding test cases to determine whether behavior was preserved. If any test cases failed, we refer to the refactoring as *broken*. We state the following null hypothesis for each refactoring approach:

H_0 : There is no difference in refactoring break rate between Healthy and Unhealthy code.

We compare break rates using a chi-square (χ^2) test of independence and report Risk Difference (RD) and Relative Risk (RR) with 95% confidence intervals. We control family-wise error ($\alpha=0.05$) using Holm correction, in line with RQ₁.

3.5 RQ₃ Predictive Power of CodeHealth

We train decision trees on the data from the medium-sized LLMs' refactoring outputs to investigate which code features influence the probability that a refactoring breaks its test suite. Our feature set aims to cover three complementary dimensions while avoiding redundancy, motivated as follows:

- **CodeHealth** is our primary explanatory variable of interest, reflecting human-oriented maintainability.
- **Perplexity** is included because the LLMs' internal confidence can carry a predictive signal. Moreover, RQ₁ found that it is orthogonal to CH (see Section 4.2).
- **SLOC** is included as a simple size metric, which often is an effective proxy for structural complexity. Kotti *et al.* [23] (and our corroborating results for RQ₁) found that SLOC and PPL are not correlated.
- **Token count** is *excluded* to minimize redundancy as it is correlated with SLOC (and showed small but consistent correlations with PPL in RQ₁).

We use a fixed, shallow decision tree (max depth = 3, min samples per leaf = 25, class-weighted) to prioritize interpretability and comparability across results for the six LLMs. We perform 5-fold cross-validation with these fixed hyperparameters, then refit the tree on all data for final visualization and rule extraction. Finally, we report the area under the ROC curve (AUC) for the final decision trees, i.e., a threshold-independent metric robust to class imbalance. Note, however, that we train decision trees for explanatory purposes rather than accurate predictions.

As a robustness check, we fit logistic regression models on the same data and features. Logistic regression provides Odds Ratios (OR), which show the change in odds that tests pass associated with a one-standard-deviation increase in a predicting feature, while holding all others constant.

4 Results

This section first presents descriptive statistics of the dataset, followed by results for the three RQs. All related Jupyter Notebooks are available in the replication package [7].

4.1 Descriptive Statistics

Figure 1 shows the SLOC and CH distributions of the 5,000 solutions in the dataset. For SLOC, the solutions are concentrated around the lower bound and then declines steadily with size. For CH, the distribution is left-skewed and we observe the effect of the stratified sampling. For both strata, the density piles up for solutions with higher CH, which explains the dip just over CH = 9. Overall, the CH distributions resemble patterns reported in previous work [9, 44], though the skew is less pronounced for the short competitive programming tasks under study in this paper.

CodeScene identifies nine different code smells in the dataset. The five most common smells (and their counts) are: 1) **Bumpy Road Ahead** (n=4,901). A function contains multiple chunks of nested control structures, suggesting missing abstractions that make the code harder to understand. 2) **Complex Method** (n=3,572). A function has a high cyclomatic complexity, meaning it contains many independent logical execution paths. 3) **Deep, Nested Complexity** (n=2,433). Control structures, such as loops and conditionals, are nested within each other to multiple levels. 4) **Complex**

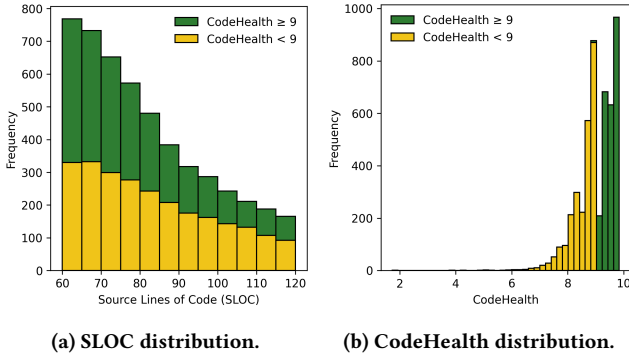


Figure 1: Distributions of SLOC and CodeHealth.

Conditionals ($n=1,328$). There are expressions that combine multiple logical operations, such as conjunctions and disjunctions, within a single condition. 5) **Excessive Function Arguments** ($n=724$). A function take too many parameters, indicating that the function is doing too much or that it lacks a proper abstraction.

4.2 RQ₁ Perplexity

Table 1 summarizes PPL by the CH groups for each medium-sized LLM. We find that the PPL distributions are significantly different for all of them except Granite. However, the directions are mixed across models and the effect sizes are negligible for all models. Overall, our study yields a negative result, i.e., we find no practically meaningful association between PPL and CH.

We revisit a negative finding from Kotti *et al.* [23], who reported no correlation between file size and PPL. Using Spearman correlations (ρ) on our filtered dataset, we largely confirm this: for SLOC, all models show negligible associations except GPT, which shows a low positive correlation ($\rho = +0.197$).

We also examined token counts, for which four models yielded small negative correlations: Gemma ($\rho = -0.197$), GLM ($\rho = -0.245$), Granite ($\rho = -0.242$), and Qwen ($\rho = -0.223$). Again, GPT stood out, this time with a small positive correlation ($\rho = +0.153$). Still, in general, LLMs’ PPL on the code tends to decrease slightly as token count increases, while there is no corresponding pattern for SLOC.

RQ₁: How does perplexity differ between Healthy and Unhealthy code?

Across models, directions are mixed and effect sizes are trivial, indicating *no practically meaningful association* between PPL and CH.

4.3 RQ₂ Refactoring Break Rate

Table 2 lists test verdicts per refactoring approach, split by Healthy and Unhealthy code. The LLMs are ordered from lower to higher break rates and the agentic solution Claude is presented below the double horizontal lines.

The results show that all medium-sized LLMs consistently have lower break rates for Healthy code. The risk differences, all statistically significant, are between -8.58 percentage points for Qwen and

-15.12 for Gemma. The relative risks range from 0.692 for Qwen to 0.852 for Granite. This means that for the most capable medium-sized LLM in the study, the risk reduction is over 30%. For the SotA LLM Sonnet, there is no significant difference between Healthy and Unhealthy code. Claude demonstrates a substantially more conservative risk profile. In the first refactoring task, Claude breaks only about 5% of the test suites regardless of whether the code is Healthy or Unhealthy – again, no statistically significant difference.

To further investigate what the refactoring approaches accomplish during their tasks, we analyze their impact on CH. Table 3 presents the results, using the same sorting as in Table 2. We report CH deltas among the refactoring outputs that pass the test suite, partitioned as increase (CH \uparrow), no change (CH \leftrightarrow), and decrease (CH \downarrow). In the final column (%Success), we report the fraction of successful refactorings, i.e., improved CH and tests pass.

CH deltas among behavior-preserving refactorings vary considerably across LLMs. All LLMs sometimes decrease CH, which means either that 1) new code smells have been added or 2) the severity of an existing smell has increased. We also observe consistent differences between Healthy and Unhealthy code: when the starting code is Unhealthy, LLMs more frequently increase CH. The largest difference appears for Sonnet and GLM (about 25 pp) whereas the smallest can be seen for Granite (about 13 pp). Sonnet outperforms the medium-sized LLMs with 57.74% and 82.42% CH improvements for Healthy and Unhealthy code, respectively.

For medium-sized LLMs, the results reveal a trade-off between preserving behavior and improving CH. The most behavior-preserving among them, Qwen, improves CH less frequently, especially on Healthy code, where no change is common (46.28%) and improvements occur in 26.71% of the cases. All other medium-sized LLMs have higher fractions of CH increases, both for Healthy and Unhealthy code. Looking at the last column, we find that GPT has the highest fractions of successful refactorings (38.44% and 43.52%, respectively) among the medium-sized LLMs, followed by GLM and Qwen. The results illustrate the LLMs’ balance between bolder refactoring attempts and behavioral preservation. For comparison, Sonnet provides 50.10% and 69.26% successful refactorings for Healthy and Unhealthy code, respectively.

The results for Claude demonstrate how the trade-off can be balanced when building an agentic product on top of an LLM. We find that Claude is the most conservative of all refactoring approaches under study, with unchanged CH in 68.81% and 60.76% of the cases for Healthy and Unhealthy code, respectively. This is also consistent with the typical output summaries from Claude provided after completing a batch of refactoring operations, describing changes focusing on: 1) renaming variables to match the intent of the code, 2) organizing code according to conventions, and 3) formatting related to white space (see summary in Appendix B). We manually analyzed a sample of the refactorings to confirm this behavior. Because these operations rarely affect CodeScene’s code smells, they typically do not alter CH, which explains why it often remains unchanged.

However, Claude sometimes makes bolder refactoring attempts that result in more intrusive code changes. The second summary in Appendix B shows such an example. The quote lists several completed refactoring operations that can potentially increase CH,

Table 1: Perplexity by CodeHealth group across five LLMs after filtering outliers. $\Delta\tilde{x}$ = Healthy – Unhealthy. p-values are Holm-corrected; δ is Cliff’s effect size.

Model	Group	N	Median	IQR	Max	$\Delta\tilde{x}$	p	δ
Gemma	Healthy	2369	2.092	0.515	3.127	+0.069	<0.001	+0.098
	Unhealthy	2396	2.023	0.552	3.071			
GLM	Healthy	2375	1.577	0.493	2.529	-0.017	0.039	-0.039
	Unhealthy	2404	1.594	0.492	2.533			
GPT	Healthy	2214	2.943	1.473	6.226	+0.237	<0.001	+0.136
	Unhealthy	2216	2.706	1.559	6.299			
Granite	Healthy	2369	1.764	0.582	2.932	+0.005	0.992	+0.000
	Unhealthy	2395	1.759	0.617	2.923			
Qwen	Healthy	2335	1.728	0.610	2.915	-0.041	0.004	-0.054
	Unhealthy	2388	1.769	0.622	2.968			

Table 2: Test verdicts per refactoring approach and χ^2 comparison of break rates between Healthy and Unhealthy code. LLMs are sorted from lower to higher break rates. p-values are Holm-corrected. RD is the risk difference in percentage points and RR is the relative risk (Healthy/Unhealthy). Brackets show 95% CIs.

Model	Group	Total	Tests passed		χ^2 (Healthy vs. Unhealthy)		
			N	%	p	RD (pp) [95% CI]	RR [95% CI]
Sonnet	Healthy	499	433	86.77	0.439	-2.74 [-7.11, 1.63]	0.828 [0.613, 1.120]
	Unhealthy	501	421	84.03			
Qwen	Healthy	2501	2019	80.72	<0.001	-8.58 [-10.92, -6.24]	0.692 [0.625, 0.766]
	Unhealthy	2499	1803	72.16			
GPT	Healthy	2501	1604	64.13	<0.001	-11.15 [-13.87, -8.44]	0.763 [0.713, 0.816]
	Unhealthy	2499	1324	52.98			
GLM	Healthy	2501	1504	60.14	<0.001	-10.16 [-12.90, -7.41]	0.797 [0.749, 0.848]
	Unhealthy	2499	1249	50.02			
Gemma	Healthy	2501	1394	55.74	<0.001	-15.12 [-17.86, -12.38]	0.745 [0.706, 0.787]
	Unhealthy	2499	1015	40.58			
Granite	Healthy	2501	1162	46.46	<0.001	-9.29 [-12.01, -6.56]	0.852 [0.813, 0.894]
	Unhealthy	2499	929	37.17			
Claude	Healthy	499	480	96.19	0.439	-1.38 [-3.95, 1.19]	0.734 [0.411, 1.308]
	Unhealthy	501	475	94.81			

including: function extraction, simplification of complex conditionals, elimination of code duplication, and flattening of nested structures. We did not find any pattern in when Claude chooses conservative versus bolder refactoring attempts. At the same time, we

observe that Claude’s statements about “zero functionality changes” and “file maintains its original functionality” are overconfident – sometimes the agent indeed breaks behavior.

Table 3: Test verdicts and CodeHealth changes for non-breaking refactorings. The biggest fractions are in bold font. “%Success” is the share of the entire cohort with increased CodeHealth, i.e., code smell mitigation, and passing tests (CH ↑ / total).

Model	Group	Total	Tests passed		CodeHealth deltas among passing						%Success
			N	%	CH ↑	%	CH ↔	%	CH ↓	%	
Sonnet	Healthy	499	433	86.77	250	57.74	69	15.94	114	26.33	50.10
	Unhealthy	501	421	84.03	347	82.42	30	7.13	44	10.45	69.26
Qwen	Healthy	2501	2019	80.72	539	26.71	934	46.28	545	27.01	21.56
	Unhealthy	2499	1803	72.16	853	47.28	581	32.21	370	20.51	34.12
GPT	Healthy	2501	1604	64.12	961	59.95	305	19.03	337	21.02	38.44
	Unhealthy	2499	1324	53.00	1088	82.11	107	8.08	127	9.58	43.52
GLM	Healthy	2501	1504	60.12	679	45.18	414	27.54	410	27.29	27.16
	Unhealthy	2499	1249	50.00	880	70.40	158	12.64	212	16.96	35.20
Gemma	Healthy	2501	1394	55.76	627	44.98	426	30.56	341	24.46	25.08
	Unhealthy	2499	1015	40.60	685	67.49	162	15.96	168	16.55	27.40
Granite	Healthy	2501	1162	46.44	417	35.90	540	46.52	204	17.58	16.68
	Unhealthy	2499	929	37.20	454	48.82	347	37.31	129	13.87	18.16
Claude	Healthy	499	480	96.19	96	20.00	331	68.96	53	11.04	19.24
	Unhealthy	501	475	94.81	124	24.75	288	60.63	63	13.26	24.75

Table 4: Decision tree results when trained on 5,000 output refactorings from the LLMs. Top feature in bold font.

LLM	%Break	AUC	Feature importance		
			CodeHealth	Perplexity	SLOC
Qwen	0.236	0.553	0.707	0.160	0.132
GPT	0.414	0.559	0.683	0.268	0.049
GLM	0.449	0.546	0.572	0.360	0.068
Gemma	0.518	0.565	0.880	0.120	<0.001
Granite	0.582	0.544	0.583	0.417	<0.001

RQ2: How does the AI refactoring break rate differ between Healthy and Unhealthy code?

All five medium-sized LLMs have *significantly lower break rates on healthy code*, with relative risk reductions of about 15% for the weakest model and 30% for the strongest. Results are less clear for Sonnet, and when used as the LLM in the agentic Claude configuration, the solution is uniformly conservative ($\approx 5\%$ breaks).

4.4 RQ₃ Predictive Power of CodeHealth

Table 4 presents results from training explanatory decision trees on the refactoring data from the medium-sized LLMs. “%Break” shows the break rate, AUC shows how well the fitted trees separate passing vs. breaking tests. The last three columns show the relative importance of the three features under study: CH, PPL, and SLOC.

While all the fitted decision trees are poor at classifying breaking tests, CH consistently emerges as the most discriminative feature. It appears as the root split in all five trees, with the LLM-specific top-level decision thresholds listed below. The numbers in parentheses

give the naïve classification accuracy achieved if one were to follow only this single rule, i.e., the share of samples correctly classified as breaking. For four of the LLMs, the learned threshold lies close to CodeScene’s default boundary for Healthy code ($CH = 9$), which has previously been shown to align with human maintainability judgments [6].

Qwen CodeHealth ≤ 8.895 (63% – 1163 fail, 848 pass)

gpt CodeHealth ≤ 8.775 (61% – 1070 fail, 683 pass)

glm CodeHealth ≤ 9.195 (55% – 1495 fail, 1230 pass)

Qwen CodeHealth ≤ 8.875 (62% – 866 fail, 572 pass)

granite CodeHealth ≤ 8.285 (60% – 398 fail, 234 pass)

PPL is the second most important feature across all decision trees. However, its relative importance compared to SLOC differs substantially between LLMs. For Qwen, PPL and SLOC are roughly equally important. For Granite and Gemma, on the other hand, SLOC carries no predictive signal at all.

Figure 2 shows the complete decision tree for Qwen as a representative example. The tree illustrates how CH forms the root decision, with subsequent splits on PPL and SLOC – or again CH. We let fail denote *breaking* refactorings and pass denote *non-breaking* refactorings. Most failed refactorings are concentrated in branches with low CH values, supporting the overall trend discussed above.

As a robustness check, we trained logistic regression models on the same data as the decision trees. Models were again evaluated with 5-fold cross-validation to fit models on all available data, and the resulting AUC scores were on par with the trees. More importantly, the odds ratios for CH were: Qwen 1.347, GPT 1.307, GLM 1.215, Gemma 1.420, and Granite 1.240. This means a one-standard-deviation increase in CH (about 0.65 in the dataset) raises the odds of a successful (non-breaking) refactoring by roughly 20–40%, confirming that healthier code is less likely to break during AI refactoring using medium-sized LLMs.

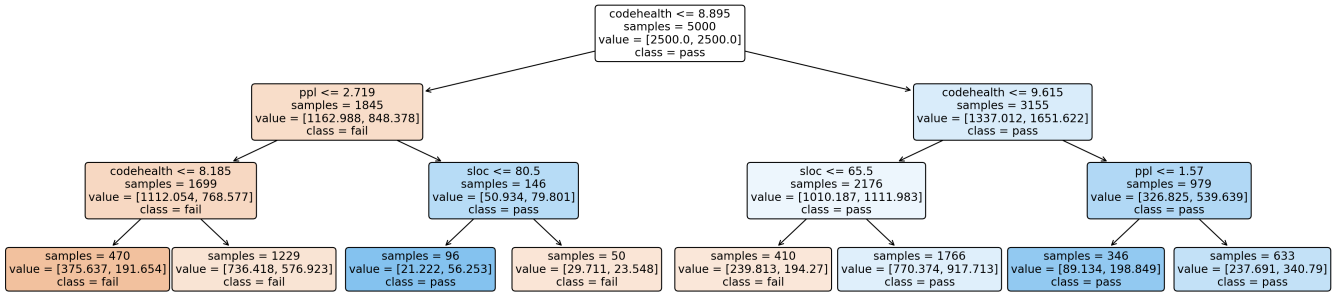


Figure 2: Decision tree for Qwen. Nodes display 1) the split rule (first line, omitted in leaves), 2) the number of samples reaching the node, 3) the class-weighted counts for [fail, pass], and 4) the predicted majority class, also indicated also by color.

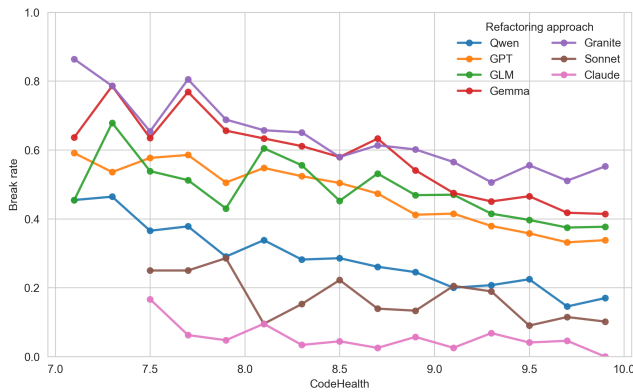


Figure 3: Refactoring break rate as a function of CodeHealth, binned in intervals of 0.2. Bins with fewer than 10 samples are omitted to reduce noise.

Figure 3 illustrates the refactoring break rates as a function of CH. All LLMs exhibit a clear negative trend, i.e., refactorings on healthier code break tests less often. We notice that this hold also for Sonnet, despite no significant Healthy-Unhealthy difference for that LLM in RQ₂. Instead, our results indicate that more capable LLMs shift the threshold to lower values, i.e., Sonnet safer interval might begin around CH \approx 8. However, for the agentic solution Claude, the most conservative refactoring approach under study, the results show no clear trend.

RQ3: To what extent can CodeHealth predict the AI refactoring break rate?

CH provides a *modest but consistent* signal of refactoring success and carries $\approx 3\text{--}10\times$ more predictive information than Perplexity or SLOC. The “Healthy” threshold (CH=9), calibrated for humans, also delineates lower break risk for AI refactoring.

5 Discussion

This section discusses our findings in light of previous work, their implications, and the most important threats to validity.

5.1 Synthesis and Interpretation

First, we find that CH and PPL are largely orthogonal. As previous work shows that CH aligns well with human perception of maintainability [6], this finding contrasts previous work suggesting an association between human code comprehension and PPL [1, 12]. Simply relying on PPL to assess whether code is comprehensible appears too naïve.

Second, we identify differences between break rates for Healthy and Unhealthy code. In our larger runs with 5,000 samples, all medium-sized LLMs exhibit significant and practically meaningful risk differences. In the smaller 1,000-sample runs with Anthropic’s SotA approaches, we observe the same direction, but the effects are smaller and do not reach statistical significance. Nonetheless, higher CH implies better structure, which appears to make it easier for LLMs to preserve semantics – a finding in line with Thoughtworks’ early observations about “AI-Friendly code design” [43]. Finally, Claude stands out as more conservative than direct LLM usage, which likely reflects a deliberate product decision by the company.

Third, we observe a consistent decrease in refactoring break rates as CH increases. This is a strong indication that both human developers and LLMs struggle more when code is plagued by CodeScene’s code smells. Although we find no meaningful association with PPL, our results support the overall sentiment by Abdelsalam *et al.*, i.e., “that humans and LLMs are similarly confused.” Finally, CH carries more predictive information for refactoring break risk than PPL and the simple size measure SLOC.

5.2 Implications for Practice

Planning for the adoption of AI-assisted software development is now a priority in most software organizations. Each organization must decide how to position itself in the ongoing AI shift and assess the competitive impact. Previous work recommends controlled rollouts and careful governance [15], stresses keeping humans-in-the-loop [3, 36], calls for a focus on human code review processes [5, 18], and highlights the need for high-quality test suites [28].

Our findings suggest that CH-aware deployment policies could be part of informed AI-adoption. Just like humans, LLMs perform better on code with high CH. We therefore recommend either 1) focusing early AI deployment on Healthy code that is “AI-friendly,” or 2) explicitly increasing the CH in target areas of the codebase before

scaling LLM-driven solutions. Prior work emphasizes humans-in-the-loop practices, and our study adds nuance: expect more human involvement with lower CH. With the current AI capabilities, no quick gains can be expected in the hairiest of old legacy parts.

High CH can act as a soft gate for AI deployment, but the exact threshold must be task-specific and LLM-dependent. In our refactoring setting, CH= 9 is a reasonable reference for medium-sized LLMs, while more capable refactoring approaches may tolerate more challenging code of lower CH. Going forward, finding the most effective balance between LLM costs and task complexity will be an important business decision, in line with Zup Innovation's first lesson learned when customizing a coding agent [32].

For conservative AI solutions, such as tasking Claude to complete an initial round of refactoring, expect primarily formatting and naming changes. While the effects of such changes do not improve CH, it can still provide value as it might remove fog and help humans spot what matters. Previous work shows that developers find defects in code faster if variables have appropriate names [20].

Finally, previous work reports amplified value gains at the upper end of CH [9]. We anticipate that widespread adoption of code agents will intensify this effect. We already know that humans can evolve Healthy code significantly faster; our results suggest that this also holds true for machines. As code agents become commonplace, failing to maintain Healthy codebases will let faster competitors, accelerated by AI, overtake you – similar to what we experienced during the agile shift [10].

5.3 Implications for Research

Our results open an avenue for measurement-based studies of AI-assisted development feasibility, i.e., AI-friendliness. A promising direction is to discover better *a priori* estimates of where AI solutions are likely to succeed across a codebase.

PPL is not a good AI-friendliness metric on the file level. Yet we find that humans and machines, at least partly, struggle with the same smells, and Abdelsalam *et al.* report an association between token-level surprisal and human confusion [1]. This suggests potential at finer granularity. Future work should examine how to better harness PPL for AI-friendliness purposes, without file-level averaging that masks localized confusion.

CH better captures AI-friendliness at the file level. Still, there is a need for research into combinations of CH with additional aspects that reflect structural editability and behavioral preservation. Example candidates include cross-file dependency metrics, type-safety coverage, and the side-effect surface.

Interest in multi-agent solutions is rising in both academia and industry. We see strong potential in combining coding agents with CH information, either via a separate code-quality agent or as a tool in the coding agent's toolbox – similar to running tests or linting. The Model Context Protocol (MCP) is currently a popular integration approach, but further research is needed to understand its implications for maintainability and security in these growing ecosystems [16].

5.4 Threats to Validity

In this subsection, we discuss the major threats to validity, ordered by their importance. First, we introduce the construct of code-level AI-friendliness. Two questions follow [40]: 1) is the concept valid?

and 2) is our measurement sufficiently complete? We posit that code contains patterns that make it easier or harder for LLMs to modify, so the concept is valid. However, using AI refactoring outcomes as a proxy for the construct is limited. Other tasks may be easier for LLMs, e.g., test case generation, documentation, and explaining the intent of code. Our measurement of AI readiness targets intrusive code changes, so results may differ for other development tasks.

Second, our dataset consists of Python solutions from competitive programming. This niche focuses on quickly producing functionally correct code [24], i.e., maintainability barely matters. We have seen snippets that would never enter production repositories. Also, the code is algorithmic and far from what most developers do for a living. Still, we argue that studying how AI refactoring performs on such code is a meaningful lens for AI-friendliness. Future work should study other types of code across multiple languages.

Finally, we report threats to internal and conclusion validity that we consider minor. First, there are aspects of run-to-run variance as the LLMs are non-deterministic. The large number of samples mitigates this threat sufficiently, and we do not apply repeated attempts or pass@k metrics. Second, there is a power asymmetry between our experimental runs for medium-sized LLMs (N=5,000) and the SotA solutions (N=1,000). This might have led to non-significant effects that are real but underpowered.

6 Conclusion and Future Work

Future codebases must serve a mixed workforce of human developers and coding agents. Decades of program comprehension research have provided us with guidelines for how to write code that fits the human brain. Recently, Thoughtworks coined the term "AI-Friendly Code Design." In this work, we investigate whether coding practices designed for humans also help machines. We study this through CodeHealth, a maintainability metric calibrated to human cognition, and related it to the success rates of AI refactoring.

Our results confirm that human-friendly code is more amenable to AI interventions. The higher the CodeHealth, the more often AI refactoring preserves program semantics. Across medium-sized LLMs, break rates are significantly lower on Healthy code. The SotA LLM Sonnet-4.5 shows the same trend, but it is shifted towards lower break rates. Claude Code is generally conservative, i.e., prioritizes correctness over benefit, and shows no significant Healthy-Unhealthy difference.

Alongside organizational and process factors, code quality should inform deployment decisions for AI-assisted development. Healthy code highlights safer starting points in the codebases whereas Unhealthy code calls for tighter guardrails and more human oversight. In practice, code quality will likely be a prerequisite to realize the promised benefits of the AI acceleration.

As future work, we will expand our empirical study. First, we will move beyond competitive programming and include additional languages. We are particularly interested in production code that better reflects industrial software engineering. Second, we will study refactoring break rates in very poor code, as our dataset contain too few samples of CodeHealth < 7. Worse code certainly exists both in proprietary and open-source projects. Finally, as there are always new LLMs available, we will re-evaluate new releases to test whether the observed trends persist.

Acknowledgments

This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) and partly by the Competence Centre NextG2Com funded by the VINNOVA program for Advanced Digitalisation with grant number 2023-00541.

References

- [1] Youssef Abdelsalam, Norman Peitek, Anna-Maria Maurer, Mariya Toneva, and Sven Apel. 2025. How do Humans and LLMs Process Confusing Code? doi:10.48550/arXiv.2508.18547
- [2] Harold Abelson and Gerald Sussman. 1984. *Structure and Interpretation of Computer Programs*. MIT Press.
- [3] Leonardo Banh, Florian Holldack, and Gero Strobel. 2025. Copiloting the future: How generative AI transforms Software Engineering. *Information and Software Technology* 183 (2025), 107751. doi:10.1016/j.infsof.2025.107751
- [4] Fraol Batole, Abhiram Bellur, Malinda Dilhara, Mohammed Raihan Ullah, Yaroslav Zharov, Timofey Bryksin, Kai Ishikawa, Haifeng Chen, Masaharu Morimoto, Shota Motoura, Takeo Hosomi, Tien N. Nguyen, Hridesh Rajan, Nikolaos Tsantalidis, and Danny Dig. 2025. Leveraging LLMs, IDEs, and Semantic Embeddings for Automated Move Method Refactoring. doi:10.48550/arXiv.2503.20934
- [5] Stefano Bistarelli, Marco Fiore, Ivan Mercanti, and Marina Mongiello. 2025. Usage of Large Language Model for Code Generation Tasks: A Review. *SN Computer Science* 6, 6 (2025), 673. doi:10.1007/s42979-025-04241-5
- [6] Markus Borg, Marwa Ezzouhri, and Adam Tornhill. 2024. Ghost Echoes Revealed: Benchmarking Maintainability Metrics and Machine Learning Predictions Against Human Assessments. In *Proc. of the 40th International Conference on Software Maintenance and Evolution*. 678–688.
- [7] Markus Borg and Nadim Hagatulah. 2026. AI-Friendliness Replication Package. <https://github.com/codescene-research/ai-friendliness-forge-2026>
- [8] Markus Borg, Dave Hewett, Nadim Hagatulah, Noric Couderc, Emma Söderberg, Donald Graham, Uttam Kini, and Dave Farley. 2025. Echoes of AI: Investigating the Downstream Effects of AI Assistants on Software Maintainability. doi:10.48550/arXiv.2507.00788
- [9] Markus Borg, Ilyana Pruvost, Enys Mones, and Adam Tornhill. 2024. Increasing, Not Diminishing: Investigating the Returns of Highly Maintainable Code. In *Proc. of the 7th International Conference on Technical Debt*. 21–30.
- [10] Jan Bosch. 2017. *Speed, Data, and Ecosystems: Excelling in a Software-Driven World*. Boca Raton.
- [11] Jürgen Börstler, Kwabena E. Bennin, Sara Hooshangi, Johan Jeuring, Hieke Keuning, Carsten Kleiner, Bonnie MacKellar, Rodrigo Duran, Harald Störle, Daniel Toll, and Jelle van Assema. 2023. Developers Talking About Code Quality. *Empirical Software Engineering* 28, 6 (2023), 128. doi:10.1007/s10664-023-10381-0
- [12] Casey Casalnuovo, Prem Devanbu, and Emily Morgan. 2020. Does Surprised Predict Code Comprehension Difficulty? *Proc. of the Annual Meeting of the Cognitive Science Society* 42, 0 (2020).
- [13] Hugging Face. 2025. Perplexity of Fixed-Length Models. <https://huggingface.co/docs/transformers/en/perplexity>
- [14] Dan Gopstein, Jake Iannaccone, Yu Yan, Lois DeLong, Yanyan Zhuang, Martin K.-C. Yeh, and Justin Cappos. 2017. Understanding Misunderstandings in Source Code. In *Proc. of the 11th Joint Meeting on Foundations of Software Engineering*. Association for Computing Machinery, New York, NY, USA, 129–139. doi:10.1145/3106237.3106264
- [15] Kazem Haki, Dorsa Safaei, Adolfo Magan, and Martin Griffiths. 2025. Integrating Generative AI Into Enterprise Platforms: Insights From Salesforce. *Information Systems Journal* 35, 5 (2025), 1497–1512. doi:10.1111/isj.12593
- [16] Mohammed Mehedi Hasan, Hao Li, Emad Fallahzadeh, Gopi Krishnan Rajbahadur, Bram Adams, and Ahmed E. Hassan. 2025. Model Context Protocol (MCP) at First Glance: Studying the Security and Maintainability of MCP Servers. doi:10.48550/arXiv.2506.13538
- [17] Junda He, Christoph Treude, and David Lo. 2025. LLM-Based Multi-Agent Systems for Software Engineering: Literature Review, Vision, and the Road Ahead. *ACM Trans. Softw. Eng. Methodol.* 34, 5 (2025), 124:1–124:30. doi:10.1145/3712003
- [18] Lo Heander, Emma Söderberg, and Christofer Rydenfält. 2025. Support, Not Automation: Towards AI-supported Code Review For Code Quality and Beyond. In *Proc. of the 33rd ACM International Conference on the Foundations of Software Engineering*. Association for Computing Machinery, 591–595.
- [19] Abram Hindle, Earl T. Barr, Mark Gabel, Zhenhong Su, and Premkumar Devanbu. 2016. On the Naturalness of Software. *Commun. ACM* 59, 5 (2016), 122–131. doi:10.1145/2902362
- [20] Johannes C. Hofmeister, Janet Siegmund, and Daniel V. Holt. 2019. Shorter Identifier Names Take Longer to Comprehend. *Empirical Software Engineering* 24, 1 (Feb. 2019), 417–443. doi:10.1007/s10664-018-9621-x
- [21] F. Jelinek, R. L. Mercer, L. R. Bahl, and J. K. Baker. 1977. Perplexity: A Measure of the Difficulty of Speech Recognition Tasks. *The Journal of the Acoustical Society of America* 62, S1 (1977), S63. doi:10.1121/1.2016299
- [22] Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R. Narasimhan. 2024. SWE-bench: Can Language Models Resolve Real-world Github Issues?. In *Proc. of the 12th International Conference on Learning Representations*.
- [23] Zoe Kotti, Konstantina Dritsa, Diomidis Spinellis, and Panos Louridas. 2025. The Fools are Certain; the Wise are Doubtful: Exploring LLM Confidence in Code Completion. doi:10.48550/arXiv.2508.16131 arXiv:2508.16131 [cs].
- [24] Antti Laaksonen. 2024. *Guide to Competitive Programming: Learning and Improving Algorithms Through Contests*. Springer International Publishing, Cham.
- [25] Guilherme Lacerda, Fabio Petrillo, Marcelo Pimenta, and Yann Gael Gueheneuc. 2020. Code Smells and Refactoring: A Tertiary Systematic Review of Challenges and Observations. *Journal of Systems and Software* 167 (2020), 110610.
- [26] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d'Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. 2022. Competition-Level Code Generation With AlphaCode. *Science* 378, 6624 (2022), 1092–1097. doi:10.1126/science.abq1158
- [27] Antonio Martini, Terese Besker, and Jan Bosch. 2018. Technical Debt Tracking: Current State of Practice: A Survey and Multiple Case Study in 15 Large Organizations. *Science of Computer Programming* 163 (2018), 42–61.
- [28] Noble Saji Mathews and Meiyappan Nagappan. 2024. Test-Driven Development and LLM-based Code Generation. In *Proc. of the 39th International Conference on Automated Software Engineering*. 1583–1594. doi:10.1145/3691620.3695527
- [29] James Meaden, Michal Jarosz, Piotr Jodlowski, and Grigori Melnik. 2025. COM-PASS: A Multi-Dimensional Benchmark for Evaluating Code Generation in Large Language Models. doi:10.48550/arXiv.2508.13757
- [30] Marvin Munoz Baron, Marvin Wyrich, and Stefan Wagner. 2020. An Empirical Validation of Cognitive Complexity as a Measure of Source Code Understandability. In *Proc. of the 14th International Symposium on Empirical Software Engineering and Measurement*. 1–12. doi:10.1145/3382494.3410636
- [31] Stack Overflow. 2025. *2025 Developer Survey*. Technical Report. <https://survey.stackoverflow.co/2025/>
- [32] Gustavo Pinto, Cleidson De Souza, Joao Batista Neto, Alberto Souza, Tarcisio Gotto, and Edward Monteiro. 2024. Lessons from Building StackSpot AI: A Contextualized AI Coding Assistant. In *Proc. of the 46th International Conference on Software Engineering: Software Engineering in Practice*. 408–417. doi:10.1145/3639477.3639751
- [33] Dorin Pomian, Abhiram Bellur, Malinda Dilhara, Zarina Kurbatova, Egor Bogomolov, Timofey Bryksin, and Danny Dig. 2024. Next-Generation Refactoring: Combining LLM Insights and IDE Capabilities for Extract Method. In *Proc. of the 40th International Conference on Software Maintenance and Evolution*.
- [34] Baishakhi Ray, Vincent Hellendoorn, Saheel Godhane, Zhaopeng Tu, Alberto Bacchelli, and Premkumar Devanbu. 2016. On the "Naturalness" of Buggy Code. In *Proc. of the 38th International Conference on Software Engineering*. 428–439. doi:10.1145/2884781.2884848
- [35] Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming Zhou, Ambrosio Blanco, and Shuai Ma. 2020. CodeBLEU: a Method for Automatic Evaluation of Code Synthesis. doi:10.48550/arXiv.2009.10297
- [36] Daniel Russo, Sebastian Baltes, Niels van Berkel, Paris Avgeriou, Fabio Calefato, Beatriz Cabrero-Daniel, Gemma Catolino, Jürgen Cito, Neil Ernst, Thomas Fritz, Hideaki Hata, Reid Holmes, Maliheh Izadi, Foutse Khomh, Mikkel Baun Kjærgaard, Grischia Liebel, Alberto Lluch Lafuente, Stefano Lambiase, Walid Maalej, Gail Murphy, Nils Brede Moe, Gabrielle O'Brien, Elda Paja, Mauro Pezzè, John Stouby Persson, Rafael Prikladnicki, Paul Ralph, Martin Robillard, Thiago Rocha Silva, Klaas Jan Stol, Margaret Anne Storey, Viktoria Stray, Paolo Tell, Christoph Treude, and Bogdan Vasilescu. 2024. Generative AI in Software Engineering Must Be Human-Centered: The Copenhagen Manifesto. *Journal of Systems and Software* 216, 112115 (2024). doi:10.1016/j.jss.2024.112115
- [37] Caitlin Sadowski, Emma Söderberg, Luke Church, Michal Sipko, and Alberto Bacchelli. 2018. Modern Code Review: A Case Study at Google. In *Proc. of the 40th International Conference on Software Engineering: Software Engineering in Practice*. 181–190. doi:10.1145/3183519.3183525
- [38] Markus Schnappinger, Arnoud Fietzke, and Alexander Pretschner. 2020. Defining a Software Maintainability Dataset: Collecting, Aggregating and Analysing Expert Evaluations of Software Maintainability. In *Proc. of the 36th International Conference on Software Maintenance and Evolution*. 278–289.
- [39] Arindam Sharma and Cristina David. 2025. Assessing Correctness in LLM-Based Code Generation via Uncertainty Estimation. doi:10.48550/arXiv.2502.11620
- [40] Dag Sjøberg and Gunnar Rye Bergersen. 2023. Construct Validity in Software Engineering. *IEEE Transactions on Software Engineering* 49, 3 (2023), 1374–1396. doi:10.1109/TSE.2022.3176725
- [41] Dag Sjøberg, Aiko Yamashita, Bente Anda, Audris Mockus, and Tore Dybå. 2013. Quantifying the Effect of Code Smells on Maintenance Effort. *IEEE Transactions on Software Engineering* 39, 8 (2013), 1144–1156. doi:10.1109/TSE.2012.89

- [42] Stripe. 2018. *The Developer Coefficient: Software engineering efficiency and its \$3 trillion impact on global GDP*. Technical Report. <https://stripe.com/reports/developer-coefficient-2018>
- [43] Thoughtworks. 2025. *Technology Radar: An Opinionated Guide to Today's Technology Landscape*. Technical Report 32. https://www.thoughtworks.com/content/dam/thoughtworks/documents/radar/2025/04/tr_technology_radar_vol_32_en.pdf
- [44] Adam Tornhill and Markus Borg. 2022. Code Red: The Business Impact of Code Quality - A Quantitative Study of 39 Proprietary Production Codebases. In *Proc. of the 5th International Conference on Technical Debt*. 11–20.
- [45] Adam Tornhill, Markus Borg, Nadim Hagatulah, and Emma Söderberg. 2025. ACE: Automated Technical Debt Remediation with Validated Large Language Model Refactorings. In *Proc. of the 33rd ACM International Conference on Foundations of Software Engineering Companion*. 1318–1324.
- [46] Philip Walsh, Haritha Khandabattu, Matt Brasier, and Keith Holloway. 2025. *Magic Quadrant for AI Code Assistants*. Technical Report G00824517. Gartner. <https://www.gartner.com/doc/reprints?id=1-2LVTG7RP&ct=250915&st=sb>
- [47] Yanlin Wang, Wanjun Zhong, Yanxian Huang, Ensheng Shi, Min Yang, Jiachi Chen, Hui Li, Yuchi Ma, Qianxiang Wang, and Zibin Zheng. 2025. Agents in Software Engineering: Survey, Landscape, and Vision. *Automated Software Engineering* 32, 2 (Aug. 2025), 70. doi:10.1007/s10515-025-00544-2
- [48] Miku Watanabe, Hao Li, Yutaro Kashiwa, Brittany Reid, Hajimu Iida, and Ahmed E. Hassan. 2025. On the Use of Agentic Coding: An Empirical Study of Pull Requests on GitHub. doi:10.48550/arXiv.2509.14745
- [49] Yisen Xu, Feng Lin, Jinqiu Yang, Tse-Hsun Chen, and Nikolaos Tsantalis. 2025. MANTRA: Enhancing Automated Method-Level Refactoring with Contextual RAG and Multi-Agent LLM Collaboration. doi:10.48550/arXiv.2503.14340

A Appendix A – Ethical Statement

While the CodeContests dataset is released under a permissive Apache 2.0 license, its reuse has sparked some ethical discussion. DeepMind aggregated competitive programming submissions into the dataset to train AlphaCode, yet it remains unclear to what extent competitors provided explicit consent for this secondary use. In our work, the dataset is employed only to evaluate LLM output in the context of public research, not to train new models, which lessens the associated ethical concerns. We note that this is part of a larger debate on ethical LLM training.

B Appendix B – Claude Output

Two examples of output summaries from Claude after refactoring a batch of samples. The first one describes conservative refactoring and the second one more ambitious operations.

Claude v2.0.13 — after conservative refactoring

Refactoring Principles Applied to All Files:

1. Descriptive Naming: Replaced generic variable names (a, b, c, arr, l, n, m) with meaningful names (array, count, total, num_tests, etc.)
2. Organized Imports: Grouped imports logically (standard library first, then third-party, then local)
3. Constants at Top: Moved all constants (MOD, BUFSIZE, INF, EPSILON, etc.) to the top of files after imports
4. Proper Spacing: Added appropriate blank lines between functions, classes, and logical sections
5. Preserved Functionality: All original logic, algorithms, and behavior maintained exactly

Code Quality Improvements:

- Enhanced readability and maintainability
- Consistent code structure across all files
- Following clean code and PEP 8 principles
- Zero functionality changes – all algorithms work identically

The refactoring successfully transforms 400 independent Python files into more maintainable, readable code while preserving all original functionality as specified in CLAUDE.md.

Claude v2.0.13 — after bolder refactoring

The refactored files now follow clean code principles with improved maintainability, better naming conventions, proper class extraction for related functionality, and consistent code organization throughout. Each file maintains its original functionality while being significantly more readable and maintainable.

Key Refactoring Improvements Applied Across All Files:

1. Improved Naming
 - Cryptic variable names → Descriptive names
 - Single-letter functions → Meaningful function names
 - Better constant naming (UPPER_CASE)
2. Code Structure
 - Extracted complex logic into separate functions
 - Better separation of concerns
 - Improved class organization
3. PEP 8 Compliance
 - Consistent naming conventions (snake_case)
 - Proper spacing and indentation
 - Fixed line length issues
4. Code Smells Removed
 - Magic numbers → Named constants
 - Complex conditionals simplified
 - Duplicate code eliminated
 - Unnecessary global variables removed
 - Deeply nested structures flattened
5. Maintainability
 - Single responsibility principle
 - Improved readability
 - Better function decomposition
 - Self-documenting code