

Refactoring vs Refuctoring

Code quality in the AI age

Over the coming decades, we'll have a hybrid of code written by both humans and machines. Who has the overall mental model in that context, and how do we ensure our AI generates human-readable code? To face the challenge, we need a safety net to enforce healthy code.

Adam Tornhill

Quality matters

Code Health

Source code

```
let inst;
if (isClass) {
  inst = new Component (element.props, publicContext, u

if (typeof Component.getDerivedStateFromProps === 'f
if (___DEV___) {
  if (inst.state === null || inst.state === undefined) {
    const componentName = getComponentName (Cor
    if (!didWarnAboutUninitializedState {componentName
      warningWithoutStack(
        false,
        '%s' uses 'getDerivedStateFromProps' but its initia
        '%s'. This is not recommended. Instead, define the
        'assigning an object to 'this.state' in the construc
        'This ensures that 'getDerivedStateFromProps' arg
        componentName,
      );
      didWarnAboutUninitializedState | componentName
    }
  }
}
let partialState = Component.getDerivedStateEmPop
null,
element.props,
inst.state,
};
if (___DEV___) {
  if (partialState === undefined) {
    const componentName = getComponentName (Cr
```

Parser

Examples on unhealthy code

Module level issues:

- **Low Cohesion:** many responsibilities
- **Brain Class:** low cohesion + large class + at least one Brain Method
- **Lack of Modularity:** too many business aspects

Function level issues:

- **Brain Methods:** complex functions which centralize the behavior of the module
- **Copy-pasted logic:** missing abstractions, DRY violations
- **Copy-pasted logic:** lack domain language

Implementation level issues:

- **Deeply Nested Logic:** if-statements inside if-statements
- **Primitive Obsession:** missing a domain language
- **Complex Conditional:** hard to understand

Score,
aggregate
and
categorize

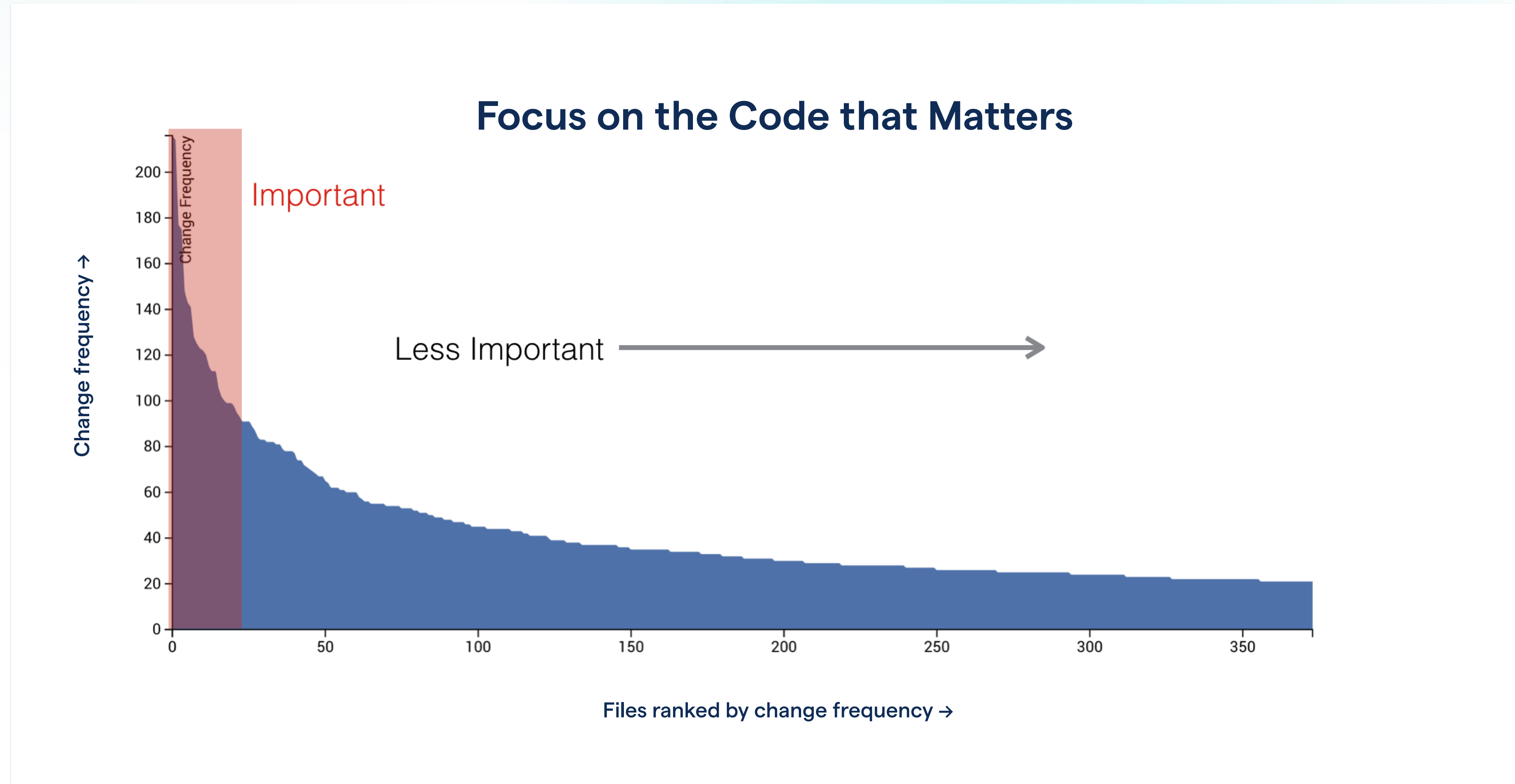
Code Health categories

Healthy code with low risk

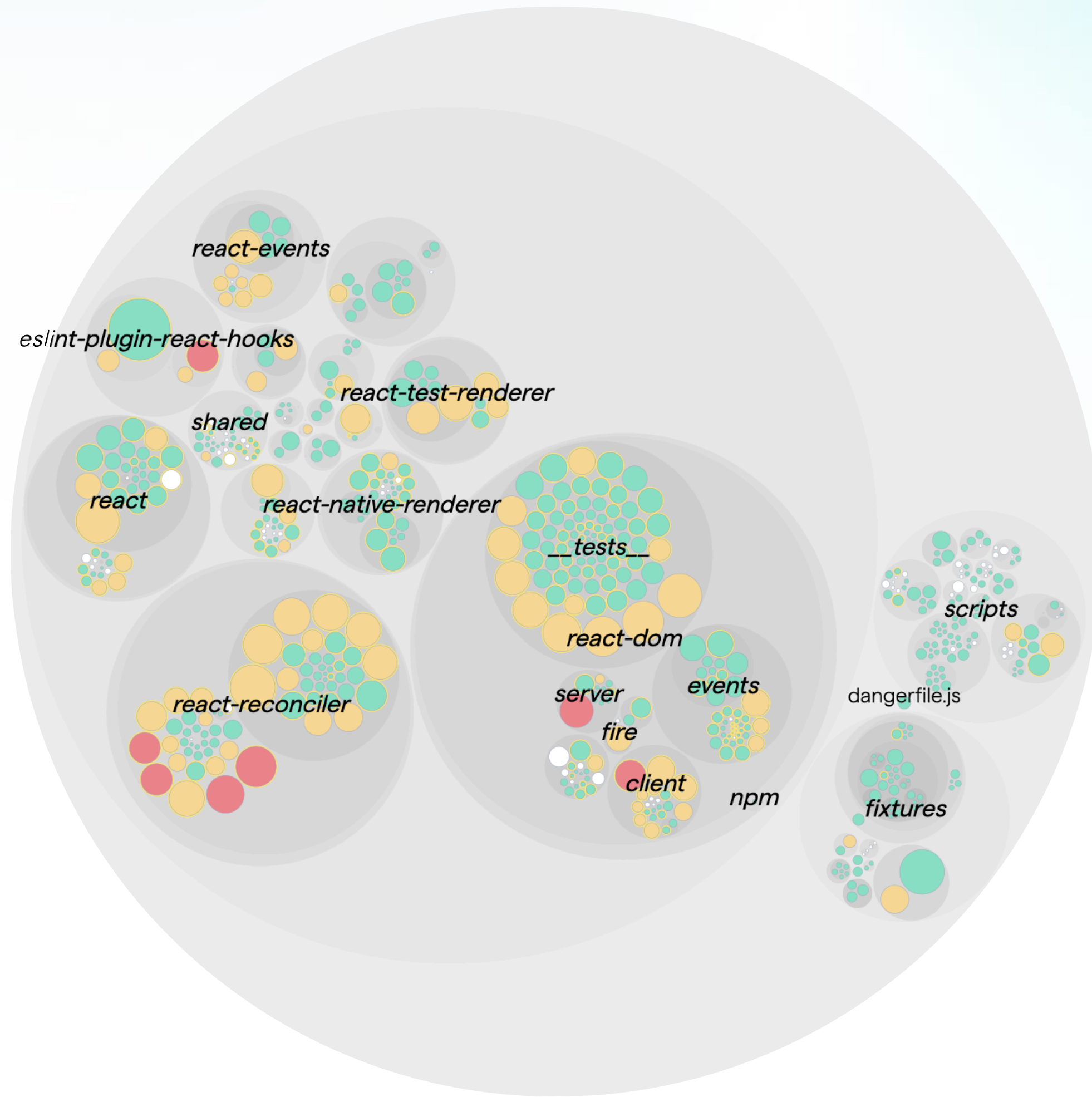
Increased maintenance efforts


Unhealthy code with significant issues and risks

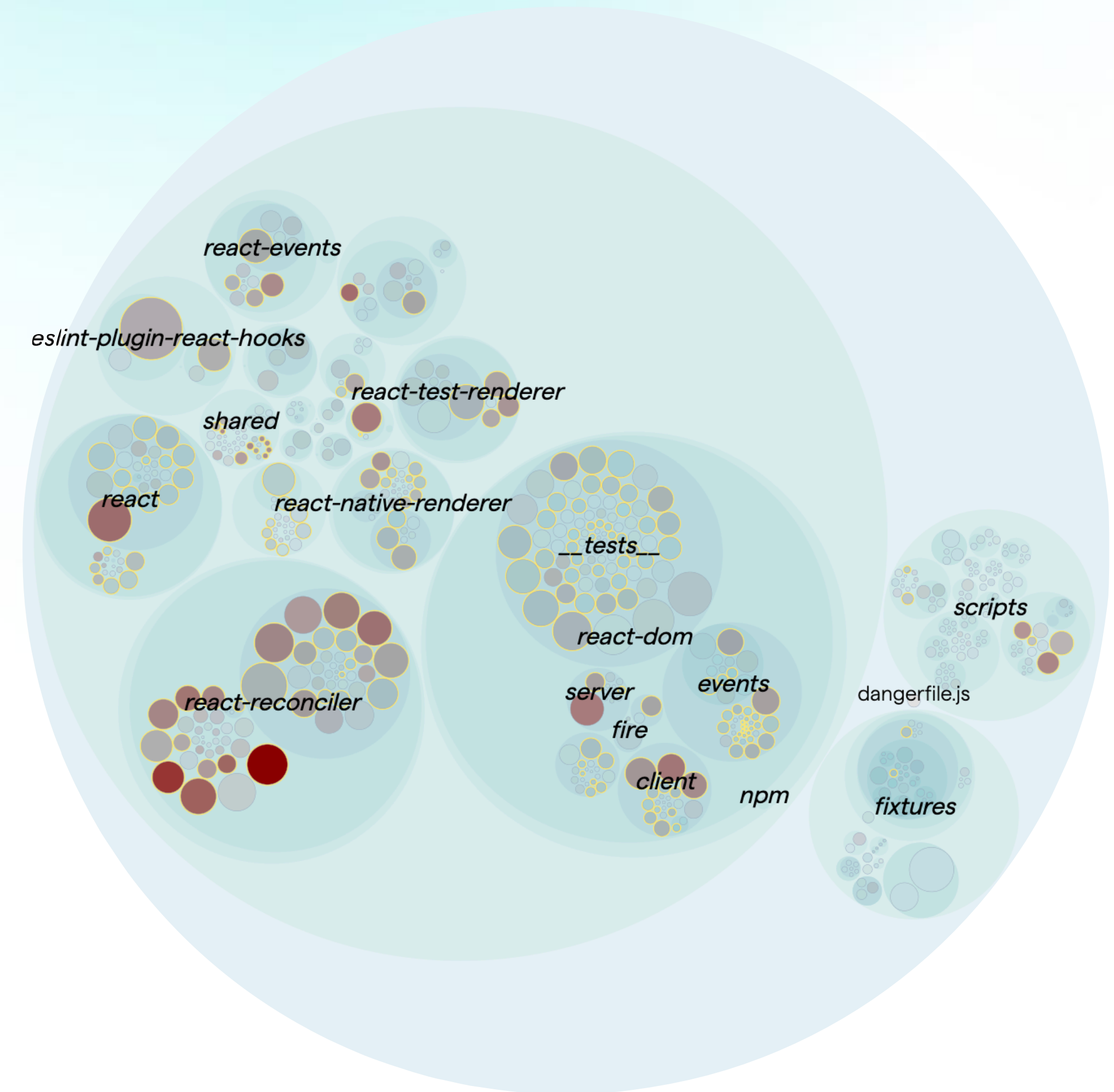
Important code

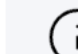
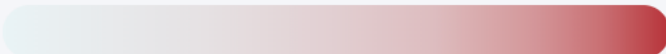


Code Health - Relevance



 Healthy Problematic Unhealthy ...



 Low development activity  Hotspot

Why all of these matter after all? - Code Red



Code Red:
The business impact of low code quality

Whitepaper

This paper presents data from a large-scale study on how code quality impacts software companies in terms of time-to-market and product experience. We conclude with an analysis of the impact and specific recommendations towards successful software development.

Target audience

- Business managers
- Product owners/managers
- Technical managers
- Tech leads
- Development teams

About CodeScene

CodeScene is the intersection of code and people, empowering companies to build great software.

CodeScene was born in 2015 when founder Adam Tornhill published the book "Your Code as a Crime Scene". It introduced a new approach to software analysis which focused on the evolution of a codebase over time.

CodeScene has become the next generation of code analysis and is used by global Fortune 100 companies in a wide variety of domains.

 CodeScene

Quantitative study of code quality impact

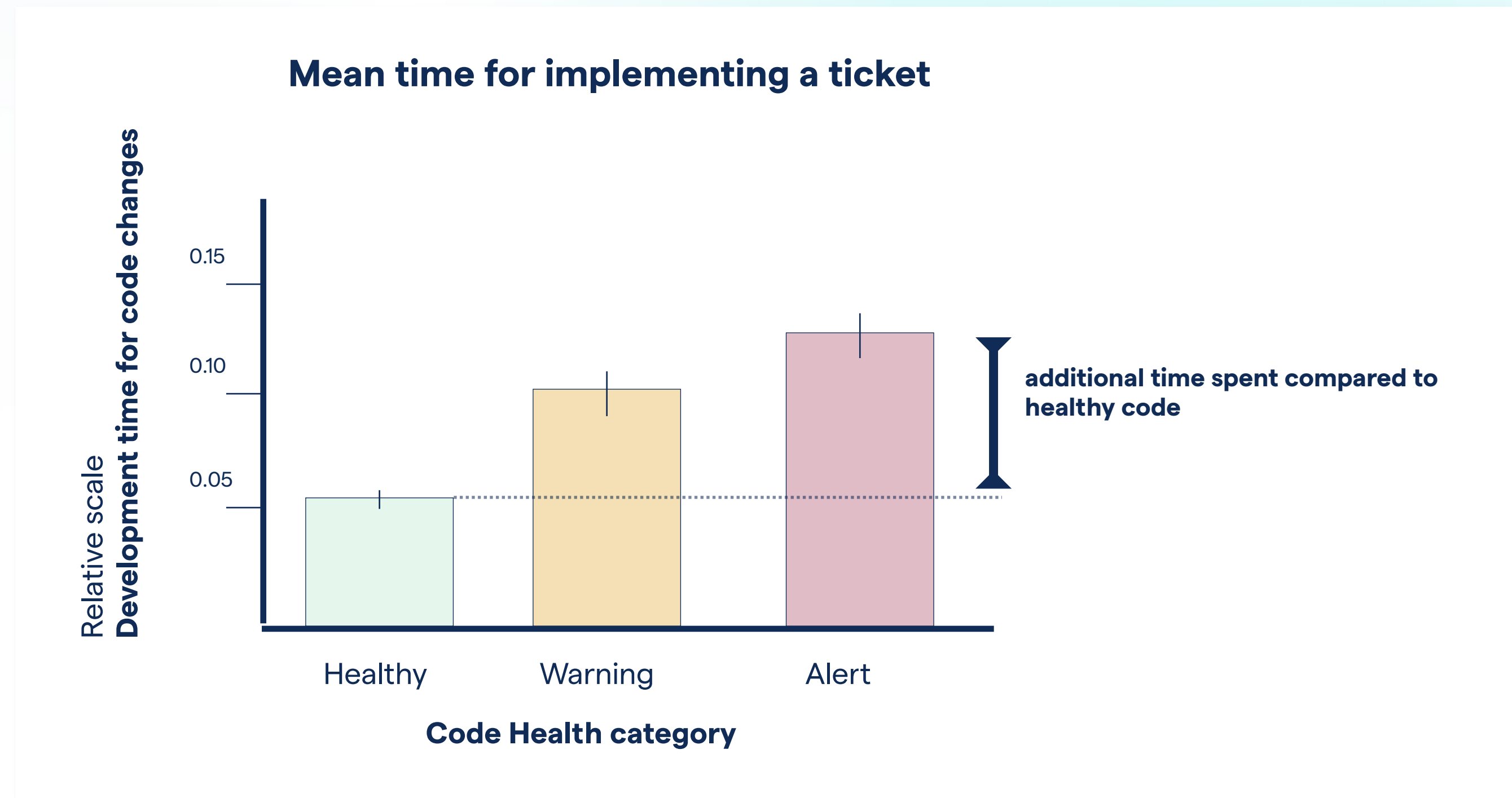
Many different industry segments

39 commercial codebases

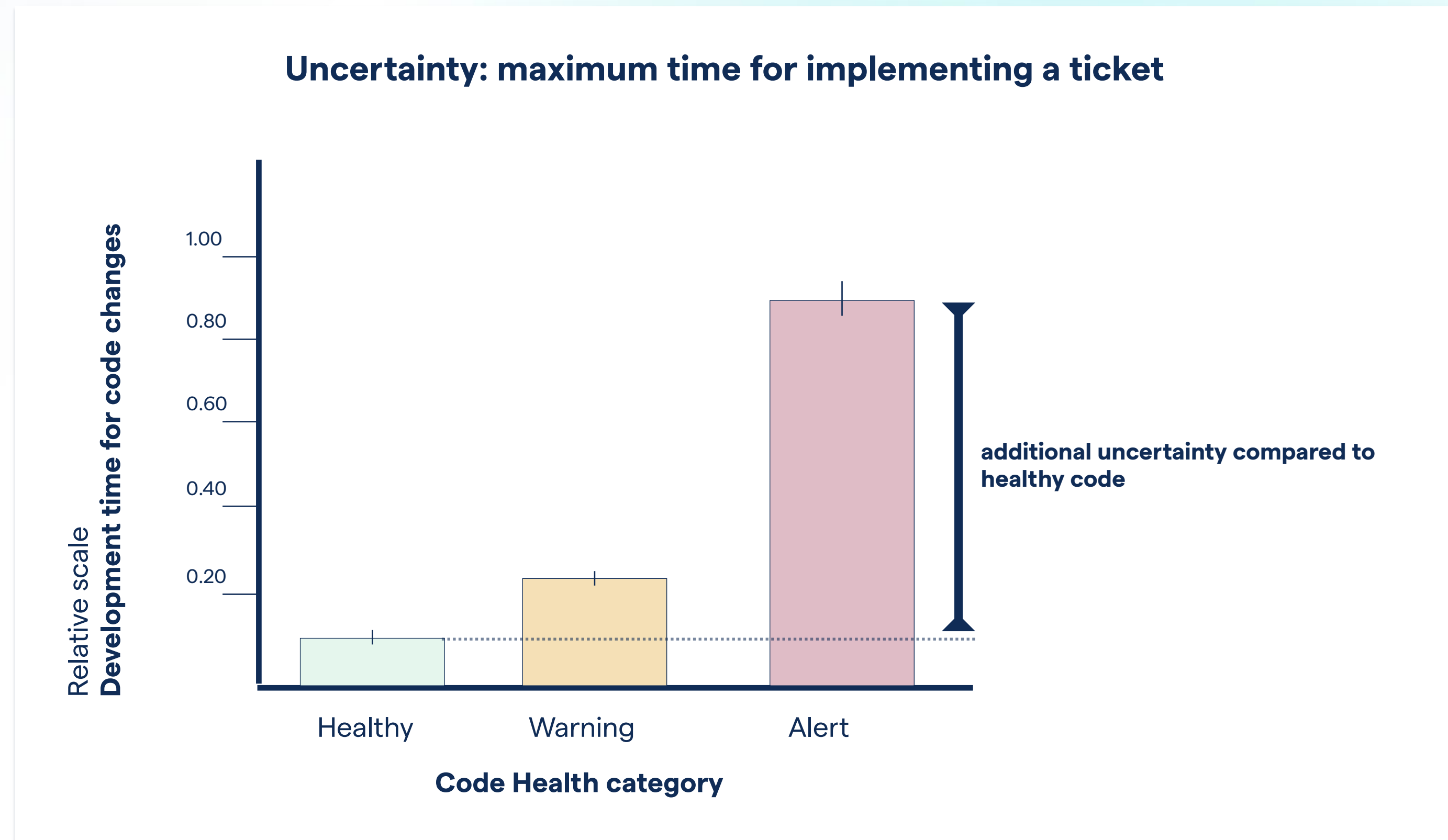
40k+ software modules

14 programming languages

Why all of these matter after all? - Code Red



Why all of these matter after all? - Code Red



Better code quality leads to faster development

Hybrid coding

The new era

Few years back, I started a new hobby...

build passing coverage 98% npm v1.24.5 docs ████████ license MIT code style standard  code health 9.45

ranjs

Requires statistically sound tests

Mathematically complex functions

The new era

Few years back, I started a new hobby...

build passing coverage 98% npm v1.24.5 docs license MIT code style standard code health 9.45

ranjs


Requires statistically sound tests

Mathematically complex functions

Code

Blame

379 lines (335 loc) · 11.5 KB

 Code 55% faster with GitHub Copilot

```
1 import { assert } from 'chai'
2 import { describe, it } from 'mocha'
3 import { repeat, trials, ksTest, chiTest, Tests } from './test-utils'
4 import { float } from '../src/core'
5 import * as dist from '../src/dist'
6 import PreComputed from '../src/dist/_pre-computed'
7 import testCases from './dist-cases'
8 import Distribution from '../src/dist/_distribution'
9
```

AI assisted coding

The Impact of AI on Developer Productivity: Evidence from GitHub Copilot

Sida Peng,^{1*} Eirini Kalliamvakou,² Peter Cihon,² Mert Demirer³

¹Microsoft Research, 14820 NE 36th St, Redmond, USA

²GitHub Inc., 88 Colin P Kelly Jr St, San Francisco, USA

³MIT Sloan School of Management, 100 Main Street Cambridge, USA

*To whom correspondence should be addressed; E-mail: sidpeng@microsoft.com.

Abstract

Generative AI tools hold promise to increase human productivity. This paper presents results from a controlled experiment with GitHub Copilot, an AI pair programmer. Recruited software developers were asked to implement an HTTP server in JavaScript as quickly as possible. The treatment group, with access to the AI pair programmer, completed the task 55.8% faster than the control group. Observed heterogenous effects show promise for AI pair programmers to help people transition into software development careers.

“Productivity benefits may vary across specific tasks and programming languages, so **more research is needed to understand how our results generalizes** to other tasks.”

AI assisted coding

The Impact of AI on Developer Productivity: Evidence from GitHub Copilot

Sida Peng,^{1*} Eirini Kalliamvakou,² Peter Cihon,² Mert Demirer³

¹Microsoft Research, 14820 NE 36th St, Redmond, USA

²GitHub Inc., 88 Colin P Kelly Jr St, San Francisco, USA

³MIT Sloan School of Management, 100 Main Street Cambridge, USA

*To whom correspondence should be addressed; E-mail: sidpeng@microsoft.com.

Abstract

Generative AI tools hold promise to increase human productivity. This paper presents results from a controlled experiment with GitHub Copilot, an AI pair programmer. Recruited software developers were asked to implement an HTTP server in JavaScript as quickly as possible. The treatment group, with access to the AI pair programmer, completed the task 55.8% faster than the control group. Observed heterogenous effects show promise for AI pair programmers to help people transition into software development careers.

“Productivity benefits may vary across specific tasks and programming languages, so **more research is needed to understand how our results generalizes** to other tasks.”

“Our results suggest that **less experienced programmers benefit more** from Copilot.”

AI assisted coding

The Impact of AI on Developer Productivity: Evidence from GitHub Copilot

Sida Peng,^{1*} Eirini Kalliamvakou,² Peter Cihon,² Mert Demirer³

¹Microsoft Research, 14820 NE 36th St, Redmond, USA

²GitHub Inc., 88 Colin P Kelly Jr St, San Francisco, USA

³MIT Sloan School of Management, 100 Main Street Cambridge, USA

*To whom correspondence should be addressed; E-mail: sidpeng@microsoft.com.

Abstract

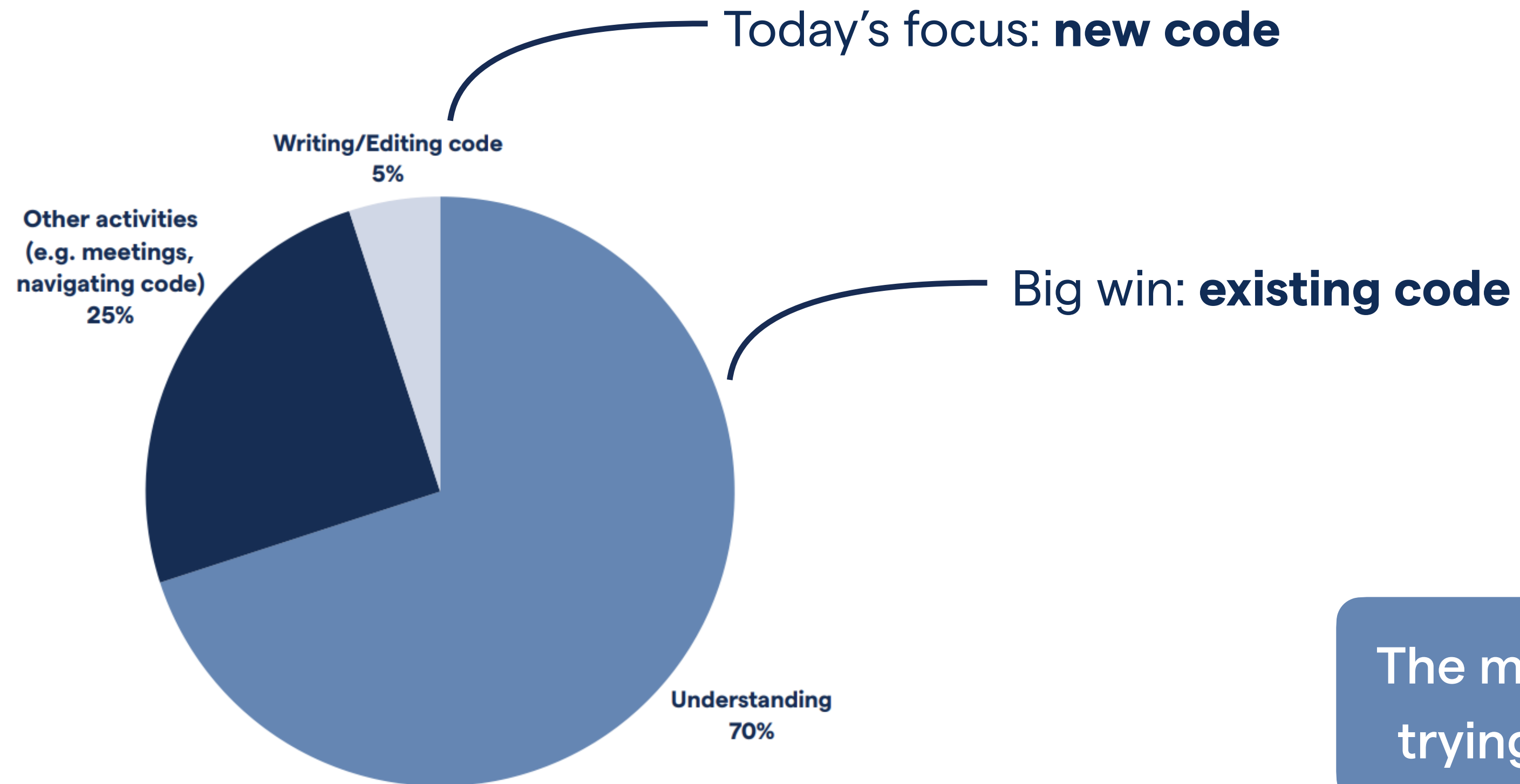
Generative AI tools hold promise to increase human productivity. This paper presents results from a controlled experiment with GitHub Copilot, an AI pair programmer. Recruited software developers were asked to implement an HTTP server in JavaScript as quickly as possible. The treatment group, with access to the AI pair programmer, completed the task 55.8% faster than the control group. Observed heterogenous effects show promise for AI pair programmers to help people transition into software development careers.

“Productivity benefits may vary across specific tasks and programming languages, so **more research is needed to understand how our results generalizes** to other tasks.”

“Our results suggest that **less experienced programmers benefit more** from Copilot.”

“Finally, this study does not examine the **effects of AI on code quality.**”

The bigger picture of time spent



55% faster on this part means ~1 hour saved per work week

The majority of a developer's time is spent trying to understand the existing system

Are we outsourcing the fun and adding to the mundane?

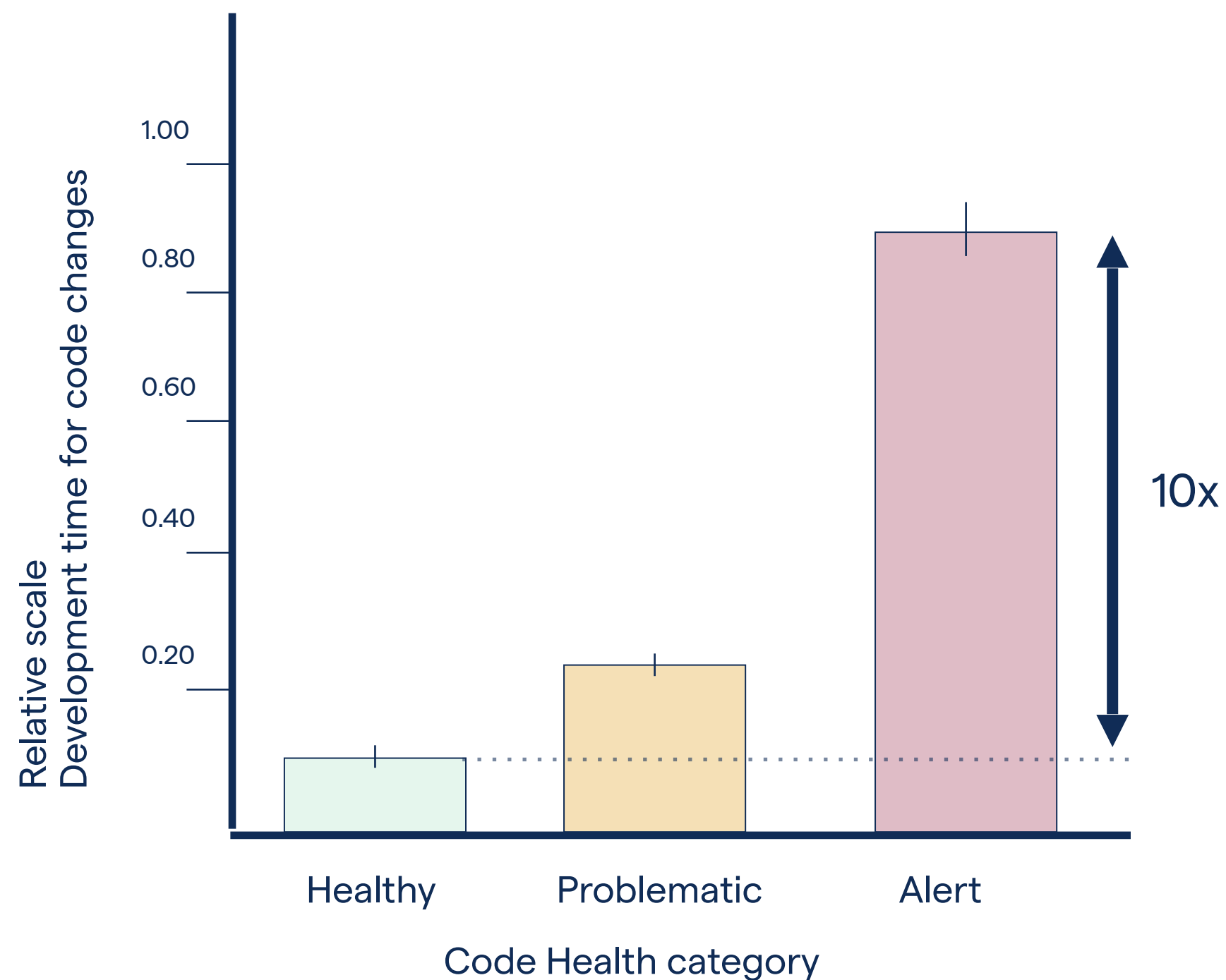
“We’ve all turned ourselves into maintenance programmers; we took the fun bit and we’re just going to give ourselves code that somebody else wrote.”

Kevlin Henney, 2024

The bigger picture

Yes, it's possible to bring that +55% to 10X

Task completions times in unhealthy code are up to 10x longer compared to green, healthy code



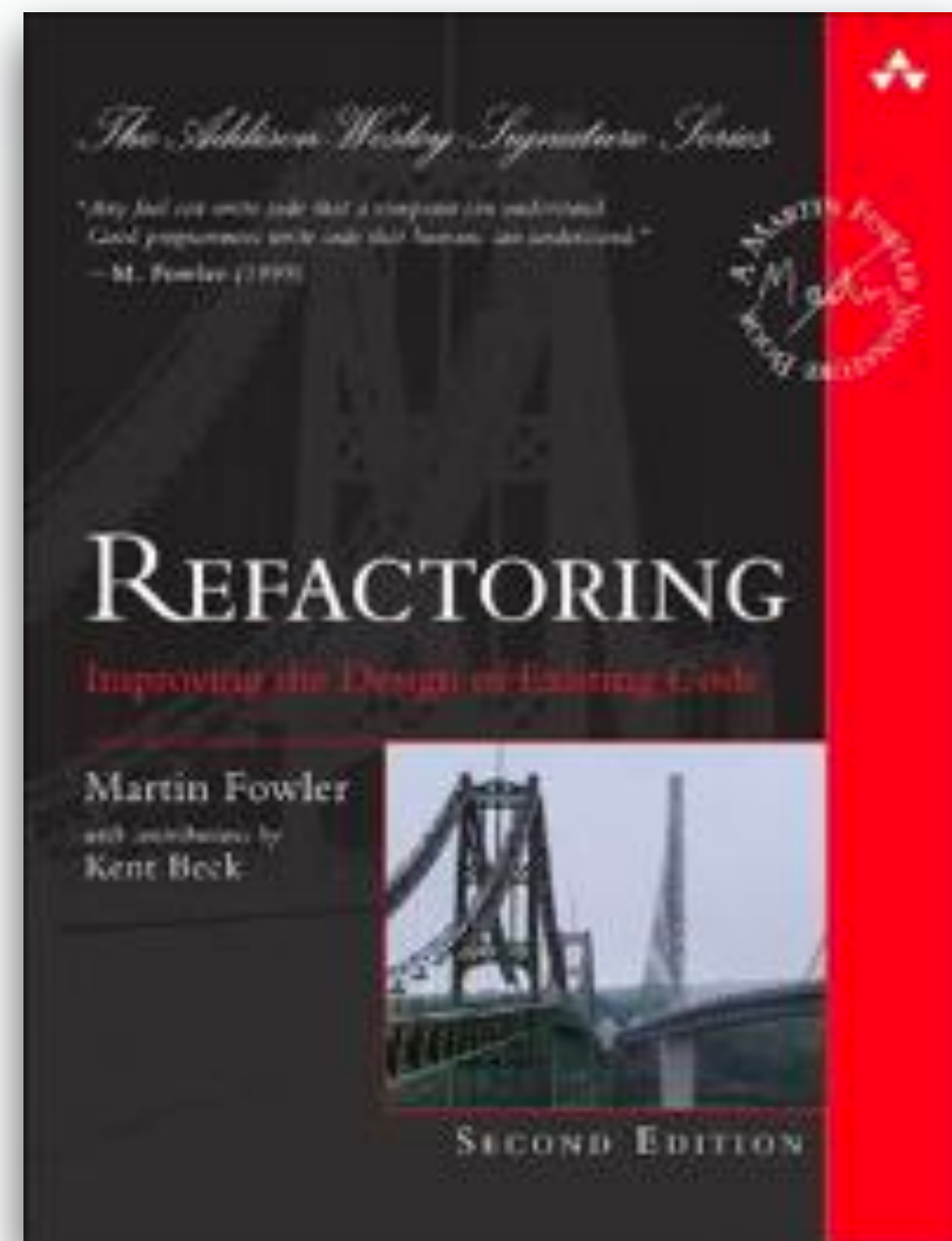
AI accelerates the creation of new code — code quality is more important than ever!

Yellow & Red Code comes with a significant on-boarding cost:

as a newcomer, you need

- **45% more time** for small tasks, and
- **93% more time** for large tasks compared to Green Code.

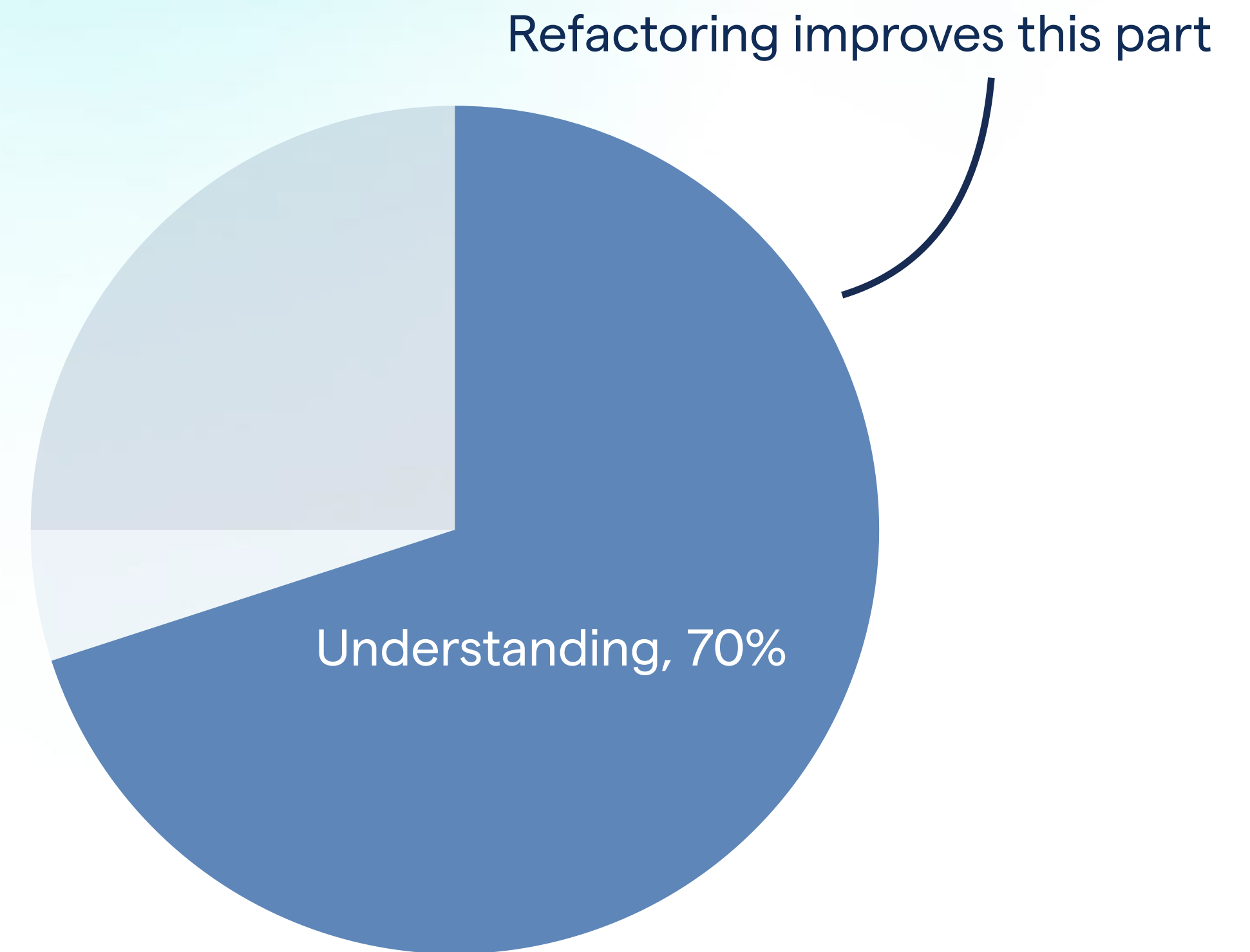
Refactoring and refactoring



Refactoring is defined as **improving the design of existing code without changing its behavior.**

- ✓ It's not a refactoring unless we improve the design.
- ✓ It's not a refactoring if we fail to preserve the behavior of the original code, e.g. we introduce a bug.

⚠ **Refactoring:** the process of changing existing code while – involuntarily – altering the program's behavior



[Research:] Let's use AI to automate refactoring

9 January 2024

Refactoring vs Refactoring:

Advancing the state of AI-automated code improvements

By Adam Tornhill, Markus Borg, PhD & Enys Mones, PhD

Summary

This report is the conclusion of a benchmark study of the most popular Large Language Models (LLMs) and their ability to generate code for refactoring tasks. We aim to illustrate the current standards and limitations, and seek to show new methodologies with higher confidence results.

100k+ refactorings generated with AI

Open source Javascript and Typescript

Benchmarking criteria: Code Health as the gold standard for code improvements

[Research:] Can AI help us improve existing code?

AI model	Valid code? (check the syntax of the refactored code)	Code Health improved? (did the code change by the AI mitigate the code smell?)	Valid refactoring? (do the tests still pass after the AI changed the code?)
PaLM 2 code [Google]	99.93%	68.75%	32.29%
GPT 3.5 [OpenAI]	100%	69.89%	30.26%
PaLM 2 text [Google]	100%	66.54%	34.73%
phind-codellama-34b-v2 [Meta, Phind]	100%	78.76%	18.14%

The average code quality

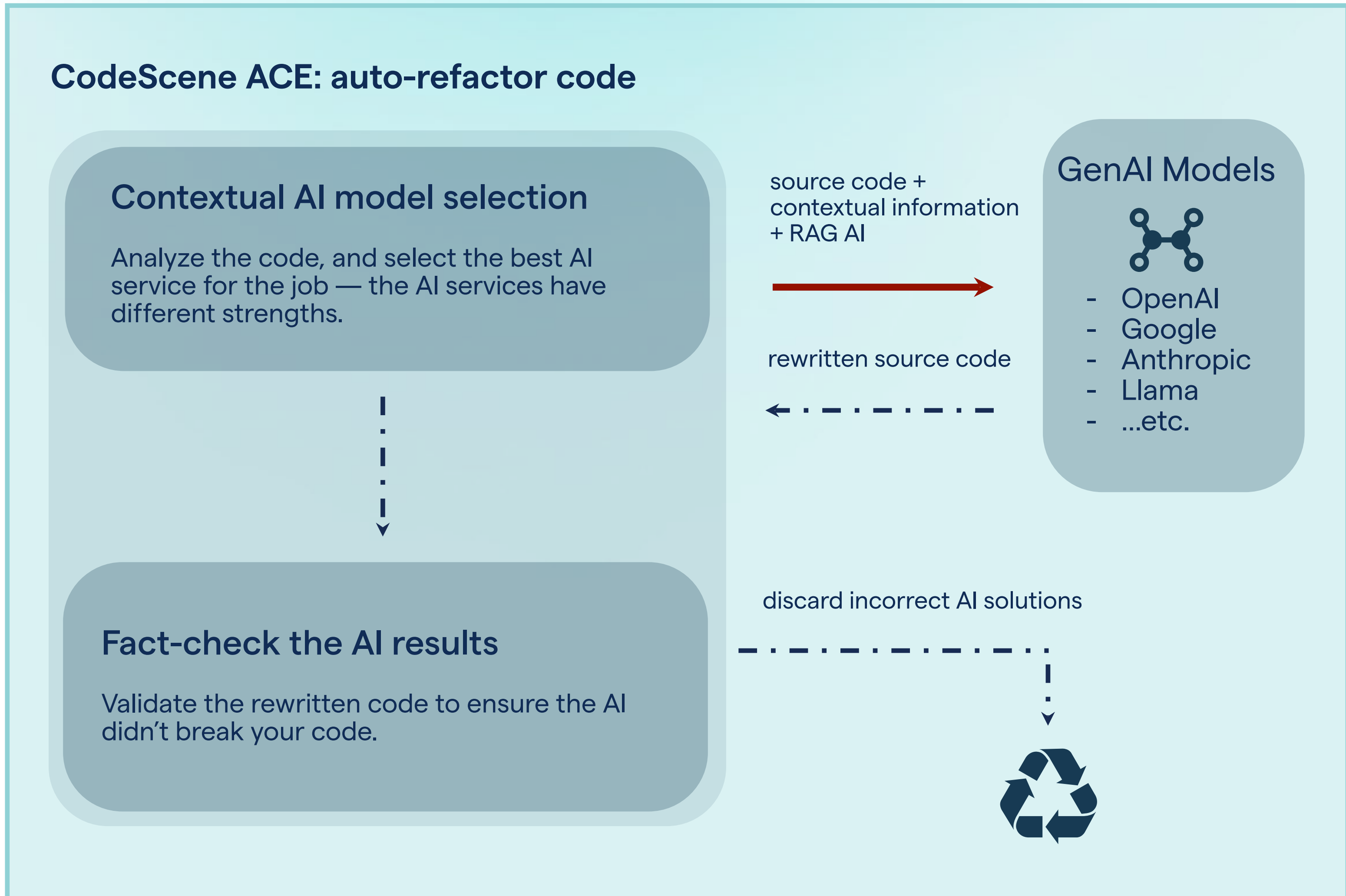
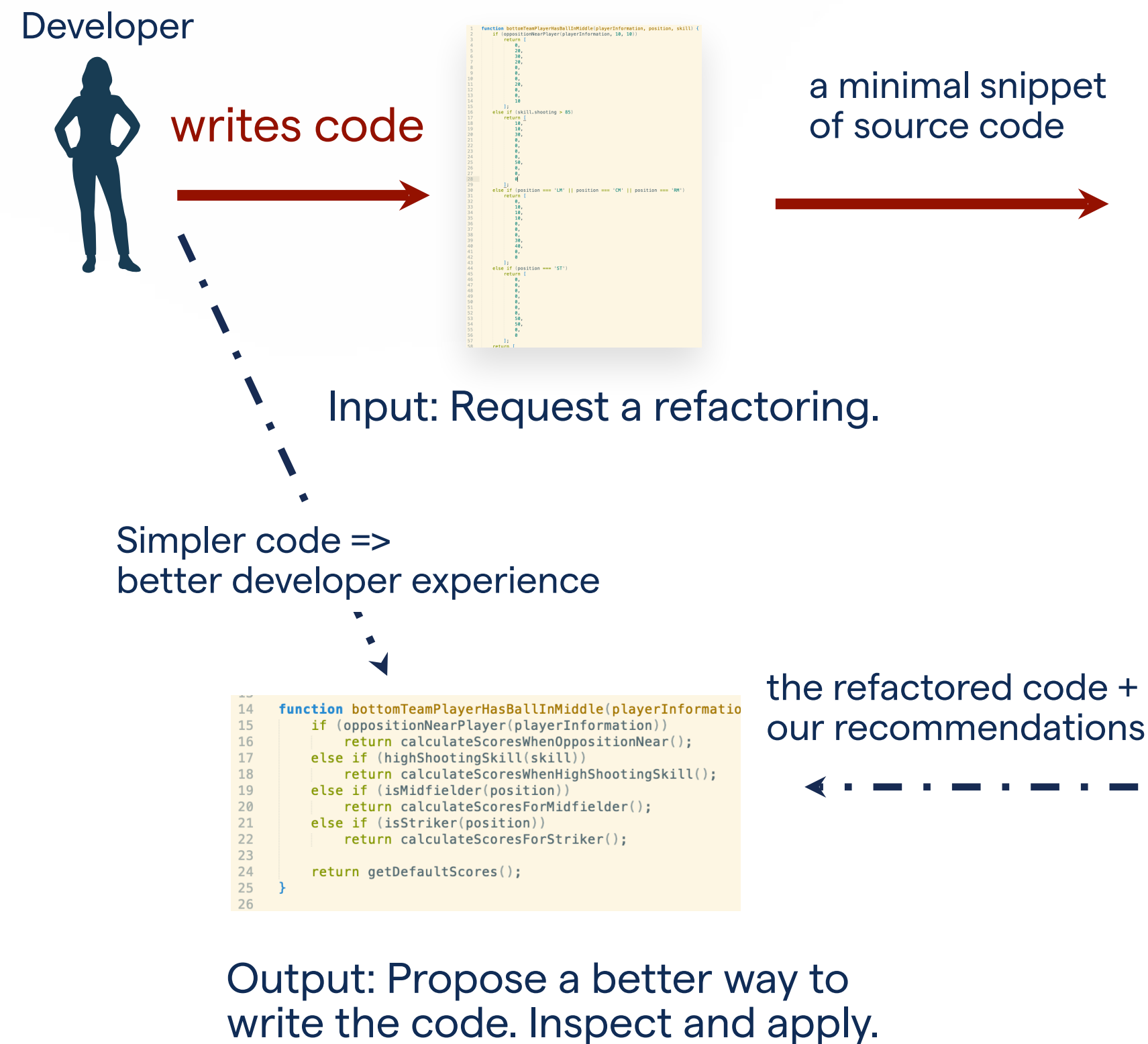
Evaluating Large Language Models Trained on Code

Mark Chen^{*1} Jerry Tworek^{*1} Heewoo Jun^{*1} Qiming Yuan^{*1} Henrique Ponde de Oliveira Pinto^{*1}
Jared Kaplan^{*2} Harri Edwards¹ Yuri Burda¹ Nicholas Joseph² Greg Brockman¹ Alex Ray¹ Raul Puri¹
Gretchen Krueger¹ Michael Petrov¹ Heidy Khlaaf³ Girish Sastry¹ Pamela Mishkin¹ Brooke Chan¹
Scott Gray¹ Nick Ryder¹ Mikhail Pavlov¹ Alethea Power¹ Lukasz Kaiser¹ Mohammad Bavarian¹
Clemens Winter¹ Philippe Tillet¹ Felipe Petroski Such¹ Dave Cummings¹ Matthias Plappert¹
Fotios Chantzis¹ Elizabeth Barnes¹ Ariel Herbert-Voss¹ William Hebgen Guss¹ Alex Nichol¹ Alex Paino¹
Nikolas Tezak¹ Jie Tang¹ Igor Babuschkin¹ Suchir Balaji¹ Shantanu Jain¹ William Saunders¹
Christopher Hesse¹ Andrew N. Carr¹ Jan Leike¹ Josh Achiam¹ Vedant Misra¹ Evan Morikawa¹
Alec Radford¹ Matthew Knight¹ Miles Brundage¹ Mira Murati¹ Katie Mayer¹ Peter Welinder¹
Bob McGrew¹ Dario Amodei² Sam McCandlish² Ilya Sutskever¹ Wojciech Zaremba¹

“We believe this is unlikely to be a large factor here, as **the GitHub dataset contains plenty of poor-quality code.**”

The bugs are designed to be of the sort we’d expect to appear commonly in the dataset; code that compiles and often runs without errors but gives an incorrect answer.”

[Innovation:] Fact-checking the AI refactorings



[Outcome:] Elevate AI to the level of human experts with a fact-checking model

	Complex Conditional	Deep, Nested Complexity	Bumpy Road	Complex Method
Raw GPT-3.5	33.7%	26.0%	26.3%	28.2%
GPT-3.5 with fact-checking	96.7%	98.4%	97.8%	98.9%

CodeScene ACE combines the results of multiple AIs and reject the incorrect solutions, **98%** of the remaining AI-generated refactorings improve the code without breaking it.

**With fact-checking, we can elevate
generative AI to achieve 10X**

Demo

How do we refactor critical code with AI and witness immediate improvement in code quality and code health?

Q&A

References

Code Red: The business impact of code quality

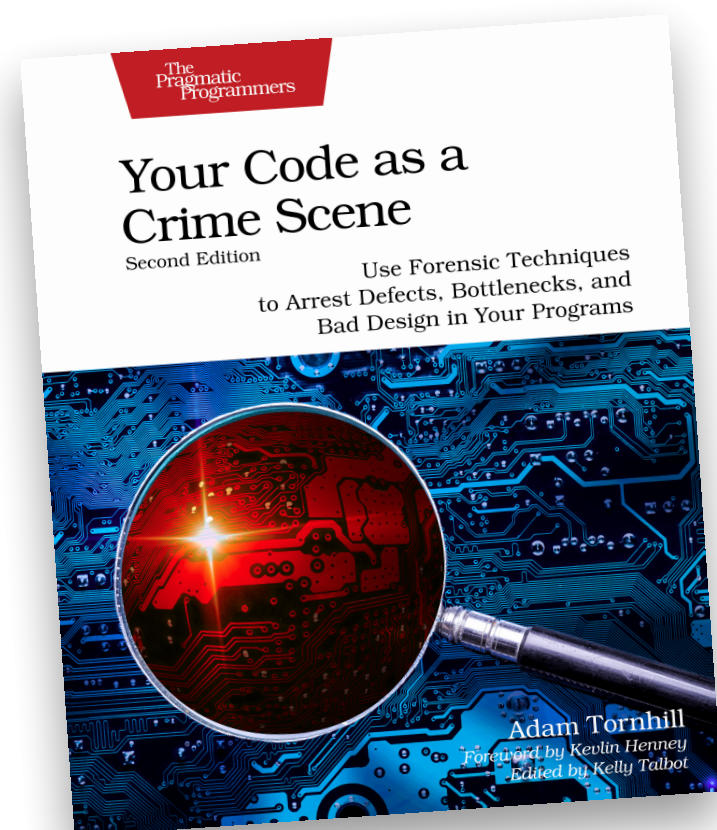
- https://codescene.com/hubfs/web_docs/Business-impact-of-low-code-quality.pdf

Refactoring vs Refactoring: Advancing the state of AI automated code improvements

- <https://codescene.com/hubfs/whitepapers/Refactoring-vs-Refactoring-Advancing-the-state-of-AI-automated-code-improvements.pdf>



[Free] Try the automated refactoring via CodeScene



Your Code as a Crime Scene, 2nd ed (2023)

<https://twitter.com/AdamTornhill>

