# Developer Perspectives on Licensing and Copyright Issues Arising from Generative AI for Coding

TREVOR STALNAKER, William & Mary, USA
NATHAN WINTERSGILL, William & Mary, USA
OSCAR CHAPARRO, William & Mary, USA
LAURA A. HEYMANN, William & Mary, USA
MASSIMILIANO DI PENTA, University of Sannio, Italy
DANIEL M GERMAN, University of Victoria, Canada
DENYS POSHYVANYK, William & Mary, USA

Generative AI (GenAI) tools have already started to transform software development practices. Despite their utility in tasks such as writing code, the use of these tools raises important legal questions and potential risks, particularly those associated with copyright law. In the midst of this uncertainty, this paper presents a study jointly conducted by software engineering and legal researchers that surveyed 574 GitHub developers who use GenAI tools for development activities. The survey and follow-up interviews probed the developers' opinions on emerging legal issues as well as their perception of copyrightability, ownership of generated code, and related considerations. We also investigate potential developer misconceptions, the impact of GenAI on developers' work, and developers' awareness of licensing/copyright risks. Qualitative and quantitative analysis showed that developers' opinions on copyright issues vary broadly and that many developers are aware of the nuances these legal questions involve. We provide: (1) a survey of 574 developers on the licensing and copyright aspects of GenAI for coding, (2) a snapshot of practitioners' views at a time when GenAI and perceptions of it are rapidly evolving, and (3) an analysis of developers' views, yielding insights and recommendations that can inform future regulatory decisions in this evolving field.

## 1 INTRODUCTION

Generative AI (GenAI) tools have become widely adopted across many different domains, including software engineering (SE) [37, 72, 106, 108]. Several GenAI coding assistants, including GitHub's Copilot [45], Tabnine [119], Codeium [24], and Cody [25], as well as general purpose tools such as ChatGPT [100], Claude [11], and Gemini [42], have become readily accessible, either as IDE extensions or standalone applications, enabling developers to perform many coding tasks with little effort, including automated code completion, summarization, and debugging. A 2024 Stack

Trevor Stalnaker, Nathan Wintersgill, Oscar Chaparro, Laura A. Heymann, Massimiliano Di Penta, Daniel M German, and Denys Poshyvanyk

Overflow survey reported that 76 percent of 65,000 participating developers were using or planning on using AI coding tools [88]; more than one million developers used Copilot [45] in its first year [129]. The popularity of these tools stems from the fact that they are easy to access and can generate diverse content with relatively little effort, thus helping users accomplish many tasks more efficiently [18, 79].

Alongside these benefits, however, using GenAI tools can introduce various issues, including bias or discrimination [38], security threats [85], and compromise of private information [54]. GenAI can also give rise to legal risk related to intellectual property concerns, including copyright infringement [43]. Many legal questions regarding using GenAI tools remain unanswered, and any answers will likely vary among jurisdictions and across cases.

As of this writing, there are several pending lawsuits against high-profile providers of GenAI tools including OpenAI [13], StabilityAI [3], and Midjourney [10], as well as active consideration by the U.S. Copyright Office and other governmental entities around the world of the copyright issues raised by the use of GenAI [31, 96, 97]. These matters implicate several considerations, including the extent to which works associated with the use of GenAI, including prompts and models, are protected by copyright and, if so, who owns the rights to such works; whether the use of open source software and other protected works as training data results in license violations or otherwise constitutes copyright infringement; and whether output resulting from the use of GenAI that is similar to preexisting work constitutes copyright infringement.

Because GenAI tools are widely used, many software developers and organizations are likely engaging in activities that implicate potential legal risks despite not having the legal training or legal advice needed to evaluate these risks. Regulation of this space– whether by governmental entities or corporate entities– should, therefore, not proceed without an awareness of the preexisting views of stakeholder groups with a vested interest in these issues. In this paper, we aim to further the discourse surrounding GenAI in the SE space and inform future decisions by identifying developers' views on (1) using GenAI tools for software development tasks and (2) associated legal concerns, focusing specifically on copyright law.

This paper presents a study conducted by a joint team of SE and legal researchers that surveyed 574 software developers worldwide who use GenAI tools for coding tasks (particularly code generation). Through an online survey and follow-up interviews, we probed the developers' opinions on potential emerging legal issues, the perception of what is copyrightable, ownership of generated code, and related considerations. We also sought to understand potential developer misconceptions, assess the impact of GenAI on their work, and evaluate their awareness of licensing/copyright risks. Using qualitative and quantitative methods to analyze survey responses, we found that developer opinions on copyright issues vary broadly, particularly on the topic of model output ownership, and that many developers are aware of the nuances and complications associated with answering these complex legal questions. We discuss the results of our study, interpret the findings under the background of U.S. law, and offer insights for future work.

In summary, we make the following contributions: (1) a comprehensive survey of 574 developers worldwide on the licensing and copyright aspects of GenAI for coding, (2) a snapshot of practitioners' views at a time when GenAI and perceptions of it are rapidly evolving, and (3) a rigorous analysis of developers' views on these topics, resulting in insights and recommendations that can help inform future policy decisions in this evolving field. We make available our survey, results, and other artifacts for transparency and validation in an online replication package [4].

## 2 BACKGROUND

Unless it is in the public domain, any textual work is generally protected by copyright in the U.S. so long as it constitutes an original work of authorship and is fixed in a tangible medium of

expression [6]. Similar provisions exist in the intellectual property legislation of other countries, including those in the European Union [27]. The owner of a copyright—which may be an individual author, joint authors, or a corporate author as the owner of a work for hire—has several exclusive rights under U.S. copyright law, including the rights to reproduce the work in copies, to create derivative works, and to distribute the work, which also includes the right to authorize others to engage in such activities [7].

The use of GenAI implicates these rights in several ways that are currently under consideration by the U.S. Copyright Office and the courts. The U.S. Copyright Office opined in February 2023— reviewing an application for copyright registration—that images generated by Midjourney were not the product of human authorship and were therefore ineligible for copyright registration [1]. The U.S. Copyright Office maintains this position [96]. There are also multiple pending cases alleging that the inclusion of copyrighted works in the training data of large language models (LLMs) and the generation of material that is substantially similar to those works constitutes infringement. For example, *New York Times Co. v. Microsoft Corporation* [121], filed in December 2023, alleges that ChatGPT uses the newspaper's content as training data, memorizes that data and then not only provides outputs that are nearly identical to that content but hallucinates content that is then incorrectly attributed to the *Times*. Another example is *Doe v. GitHub* [35], filed in November 2022, which alleges that Copilot violates the copyright licenses of the plaintiffs' open source software by using that code as training data and generating code that is a near-identical reproduction of the plaintiffs' code but without adhering to the terms of the associated licenses. Yet another lawsuit, filed by a group of authors, was brought against the companies Meta, Microsoft, and Bloomberg, alleging that the authors' works were used for training without their permission [20].

The resolution of these and similar cases depends on both technological and legal questions. On the technological side, liability may depend on whether authors can successfully show that their works are indeed included in a particular LLM's training data, as well as whether they can show that those works are copied for any length of time during the training process (or whether the training simply involves developing model parameters from the training data without reproducing it). On the legal side, liability may depend on whether and how the generated content is similar to any particular author's work. For example, generated content that is similar to previous works only with respect to the underlying concept or function is not infringing because the protections granted by U.S. copyright law for an original work of authorship do not extend to any "idea, procedure, process, system, method of operation, concept, principle, or discovery, regardless of the form in which it is described, explained, illustrated, or embodied in such work" [8].

Additionally, some legal scholars argue that the use of a work protected by copyright as part of a model's training data may be deemed a fair use under U.S. law [77], similar to how the use by Google of copyrighted material as the database for its Google Books search engine was held to be fair use by the U.S. Court of Appeals for the Second Circuit in 2015 [12]. Part of this consideration under U.S. law is likely to be whether the challenged use has a "transformative" purpose or is considered to be an unlawful derivative work; the 2021 U.S. Supreme Court case of *Google LLC v. Oracle America, Inc.* [48] and the 2023 U.S. Supreme Court case of *Andy Warhol Foundation for the Visual Arts, Inc. v. Goldsmith* [2] offer contrasting examples.

Legislation both within and outside the U.S. also has implications for GenAI and copyright. Various U.S. states have attempted to enact legislation to govern GenAI [21, 26]. In March 2024, the European Parliament approved the EU AI Act [31], which, among other things, requires the providers of general-purpose AI models to "put in place a policy to comply with Union law on copyright and related rights, and in particular to identify and comply with, including through state-of-the-art technologies, a reservation of rights expressed pursuant to Article 4(3) of Directive (EU) 2019/790" [30]. In 2023, an early announcement of the EU Global Principles for Artificial

Intelligence that led to the EU AI Act encouraged some organizations, including French media organizations France Médias Monde and TF1, to forbid companies such as OpenAI from mining their data [66]; that year also saw a temporary ban on ChatGPT in Italy [89]. Even countries like Japan, which initially took more relaxed stances, affirming that training AI models on protected content was not a violation of copyright law [74], have started to consider additional restrictions designed to protect copyright holders [113]. Gervais *et al.* [44] provide an in-depth overview of the current copyright landscape around the globe.

Prior research has explored several questions and concerns regarding copyright and GenAI [53, 83]. For example, Lee *et al.* [75] outline the intricate GenAI supply chain and its intersection at various points with aspects of U.S. copyright law, revealing that answers to copyright questions may depend on the specific details of how an individual model was trained. Craig [33] has cautioned against applying or expanding copyright principles to address issues raised by GenAI without first carefully considering the unintended consequences of such an approach. Kugler [73] provides a brief overview of the copyright challenges posed by GenAI and how various governments around the globe have begun addressing them.

Other research has considered the copyright implications of the emerging base model and subsequent fine-tuning paradigm [131] and has studied the applicability of fair use for foundation models [57]. Additional research [94] presents a review of the discourse surrounding copyright and privacy issues related to LLMs. The research most closely related to our study is by Liang *et al.* [79], who investigated why developers choose to use AI programming assistants and the challenges they encountered.

Additional work relevant to GenAI and copyright issues focuses on research that may help answer some of the above technological questions. For example, techniques exist to potentially reveal whether a machine learning model was trained on a certain data record, such as membership inference attacks [36, 60, 86, 132, 133]. Such attacks can extract training data from the target model's output, including private and personal information [22]. Memorization of training data is a known issue, including in models trained on code [128], with research suggesting that typical measures aimed at preventing verbatim extraction of training data are not sufficient to prevent the data from being obtained with other methods, such as by utilizing style-transfer prompts [63]. Tools like CopyrightCatcher [29] have been developed to detect the presence of copyrighted material in model output, evidencing that this is an issue of concern to developers.

Our work builds on this related work by contributing an in-depth exploration of developers' perceptions of the copyright issues surrounding the use of AI-based code-generation tools, both in the U.S. and elsewhere around the world. Understanding these perspectives allows future regulation in this area to be informed by beliefs on the ground rather than in a vacuum.

## 3 STUDY METHODOLOGY

The *goal* of this study is to investigate practitioners' perceptions about licensing and copyright issues related to the use of GenAI technology for software development. (Although we use the term "software development" throughout the paper, our study primarily focuses on the use of GenAI for code generation as the main task that can lead to copyright and licensing issues.) The primary *quality focus* is the mitigation of legal consequences due to copyright and licensing violations that arise when developing software with the assistance of GenAI. The *context* consists of 574 developers of open-source projects hosted on GitHub, whom we surveyed via an online survey and follow-up interviews.

The study addresses the following research questions (RQs):

• **RQ₁:** *How do developers use GenAI technologies to develop software?* Studying *how* generative AI is used by developers allows us to put the answers provided by developers in response to the subsequent RQs in context.
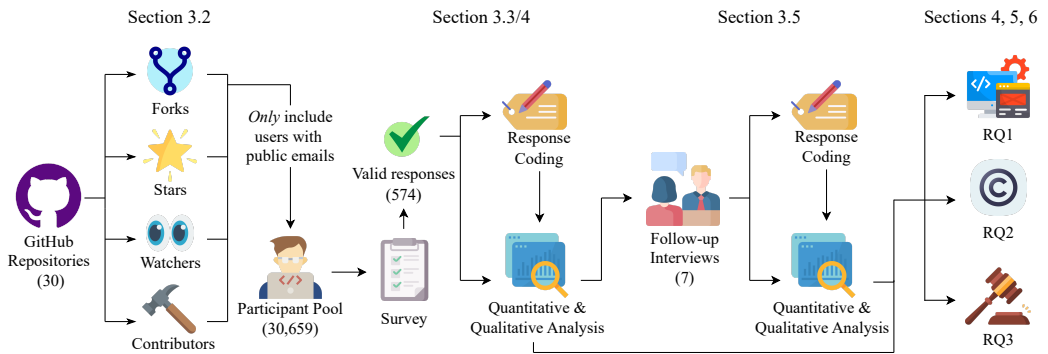
Fig. 1. Overview of the study methodology

- **RQ$_2$:** *What are developers' perceptions regarding the licensing and copyright issues that arise from the use of GenAI tools?* This RQ analyzes several dimensions of the investigated phenomenon, including developers' perceptions of potential emerging legal issues, copyrightable subject matter, and copyright ownership of generated code.

- **RQ$_3$:** *What other legal concerns do developers anticipate as the use of GenAI increases?* There are a plethora of other potential legal issues to consider as GenAI becomes more widely adopted, including data privacy, generation of malicious content, and tort liability. In this RQ, we analyze developers' perspectives on these topics.

Our study, including the survey questionnaire, the participant identification procedure, and survey and interview protocols, was approved by our institution's ethical review board. An overview of our methodology can be seen in Figure 1.

### 3.1 Survey design

To design our survey, we followed general guidelines [52] as well as SE-specific best practices [67–71, 105]. The final questionnaire went through multiple iterations of review and improvements from six SE researchers and one law researcher. We carefully formulated the questions and ensured they were written clearly and concisely to avoid biasing respondents. We also conducted a pilot study with graduate students from our research lab who actively use GenAI tools. Based on their feedback, we further refined our questionnaire, improving questions that were ambiguous or hard to understand.

The survey was structured into three core sections, as depicted in Figure 2:

- The section on "Current use of GenAI for coding" (12 questions) asked about respondents' experience using GenAI tools for code generation, including the tools they used; the benefits and challenges of using such tools; the development tasks for which they used the tools; organizational policies of using such tools; and procedures to document tool usage, among other aspects.

- The section on "Understanding and Perception of Copyright Issues" (14 questions) asked about developers' perceptions of copyright issues surrounding GenAI, including the use of their code in model training without permission; having models produce code similar to theirs; reusing without permission a prompt that they created; ownership and copyrightability of generated code; whether the use of source code in training models should require attribution and/or monetary compensation; and the effects of these concerns on developers' workflows.

- The "Demographics" section (8 questions) asked about respondents' experience with SE and AI, whether they develop open source or proprietary software, their location, and their primary programming languages, among other questions.

| Disclaimer, research procedure, participation risks, confidentiality, contact person, and research protocol info |
|---|
| **Consent** |
| **Current use of GenAI for coding (12)** <br><br> (C1) Familiarity [yes/no] <br> (C2) Use of tools [yes/no] <br>      if no: (C8) Reasons for not using GenAI for coding [open-ended] <br> (C3) Which tools [multiple choice] <br> (C4) Development tasks automated by AI-based tools [multiple choices] <br> (C5) Benefits [open-ended] <br> (C6) Problems and gaps [open-ended] <br> (C7) Whether one reads the Terms of Service [always/sometimes/just skim/never] <br> (C9) Whether the organization permits the use of GenAI for coding [yes/no/partially/don't know] <br>      if no: (C12) Reasons for not permitting use of GenAI for coding [open-ended] <br> (C10) Presence of a process to document the use of GenAI for coding [yes/no/don't know] <br>      if yes: (C11) What information is recorded [multiple choices] |
| **Understanding and Perception of Copyright Issues (14)** <br><br> (U1) Whether the use of GenAI for coding raises copyright issues [agreement 5-point Likert scale] <br> (U2) Whether litigation might limit the use of GenAI for coding [agreement 5-point Likert scale] <br> (U3) Whether the organization's use of GenAI for coding creates legal risks [agreement 5-point Likert scale] <br> (U4) The extent to which developers' workflow should account for potential legal risks [agreement 5-point Likert scale] <br> (U5) Whether copying source code should lead to monetary compensation [agreement 5-point Likert scale] <br> (U6) Whether copying source code should lead to attribution [agreement 5-point Likert scale] <br> (U7) Feeling about having one's code used as training data by the model [pleased/displeased 5-point Likert scale] <br> (U8) Feeling about having the models produce code similar to one's own code [pleased/displeased 5-point Likert scale] <br> (U9) Feeling about somebody reusing one's prompt [pleased/displeased 5-point Likert scale] <br> (U10) Comment on or expand on any answers to the questions in this section [optional, open-ended] <br> - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - <br> (U11) Opinion on who owns AI generated code [training set owner/ML model creator/prompt creator/nobody/don't know] <br>      (U12) Confidence level on the answer provided [5-point Likert scale] <br> (U13) Rationale for the interpretation of ownership of AI-generated code [open-ended] <br> (U14) Other development tasks where AI-based tools could lead to copyright/legal issues [open-ended] |
| **Demographics (8)** <br><br> (D1) Age [multiple choice with ranges] <br> (D2) Years of development experience [multiple choice with ranges] <br> (D3) Year of experience with AI-based tools [multiple choice with ranges] <br> (D4) Programming languages being used [multiple choices] <br> (D5) Nature of software being developed [proprietary, open source, educational work] <br> (D6) Country where the organization (or the person, if freelance) is based [free-text] <br> (D7) Whether the person has received training on copyright [yes/no] <br>      if yes: (D8) By whom [employer/individual/both] |
| **Whether respondent was willing to be contacted for an interview, and, if yes, contact info** |

Fig. 2. Survey design overview

The survey was designed to be completed in about 15 to 20 minutes and consisted mostly of multiple-choice questions (26), complemented by open-ended questions (7) that allowed participants to elaborate on their responses. Not every participant saw every question, as we included logic in our survey to ask only those questions relevant to a participant's indicated experience. The full survey text and a brief description are found in our replication package [4].

## 3.2 Participant identification

We sought to learn from developers who used AI tools for code generation. To this end, and consistent with prior work [28, 40, 55, 61, 64, 65, 78–80, 93, 111, 117, 127], we used GitHub as a source to identify potential participants, following multiple steps.

Table 1. Response totals

| Response Type | Count |
| --- | --- |
| All | 772 |
| Complete | 580 |
| Valid | 574 |
| Familiar with AI tools | 554 |
| Incorporated AI into workflow | 497 |

We began by searching GitHub to identify repositories related to code generation via AI. We did this by developing a list of relevant tags, including "coding-assistant" and "ai-code-generator." The full list of tags can be found in our replication package [4]. This resulted in 27 repositories, including open-source coding assistants such as "Tabby" [118] and "cptX" [32] as well as "full development" solutions such as "gpt-engineer" [50] and "gpt-pilot" [49]. We included an additional three repositories we believed would be relevant: Meta's CodeLlama, [15], GitHub Copilot [45], and Codeium [24]. This brought our total repository count to 30. Next, we sought to identify developers who had contributed to or shown interest in any of these repositories. To do this, we used the GitHub API [46] to collect account information for users that forked, starred, watched, or contributed to any of the 30 aforementioned repositories. Lastly, in order to respect the privacy of developers, we filtered this information to eliminate users who did not provide *publicly available* contact information on their GitHub profile. This resulted in 30,659 public-facing email addresses, which we used to contact individuals requesting participation in our research study. All information reported in this paper, including demographic information, was contributed voluntarily by respondents who agreed to participate in response to our request.

### 3.3 Survey response collection and analysis

Survey responses were collected using Qualtrics [5]. The survey was kept open for six weeks starting on January 21, 2024. Survey invitations were sent out via email during that period.

We obtained a total of 772 survey responses, 580 of which were complete, as shown in Table 1. We removed five responses written in languages other than English, which would have required reliance on translation tools that might not be reliable for domain-specific tasks [56]. One additional response was removed for providing irrational answers to the open-ended questions. This left us with 574 complete and valid responses.

Survey responses to five open-ended questions were analyzed through a qualitative coding approach [114]. Two SE researchers performed *open-coding* by independently assigning one or more *codes* to each response using a shared spreadsheet and codebook. Each annotator independently coded all 574 responses, adding new codes to the codebook as necessary. Once the initial coding was completed, the annotators met to settle disagreements and consolidate the set of codes. Our replication package [4] contains the final codes and definitions.

We did not base our analysis on inter-rater agreements because multiple codes could be assigned to each response, and no list of codes existed before the start of coding (*i.e.*, our approach was fundamentally inductive). However, we carefully followed best open-coding practices [114], and we leveraged coders' discussions to ensure the results' reliability. All of the responses were coded in a single iteration by the two annotators. To mitigate agreement by chance, the labels assigned to each answer were reviewed by the annotators, including those without disagreement. We avoided defining "umbrella" codes by creating complete definitions for each code and sharing them between the annotators during review, as well as by assessing the appropriateness of codes during reconciliation of the annotators' results such that codes that were too similar could be merged

Trevor Stalnaker, Nathan Wintersgill, Oscar Chaparro, Laura A. Heymann, Massimiliano Di Penta, Daniel M German, and Denys Poshyvanyk
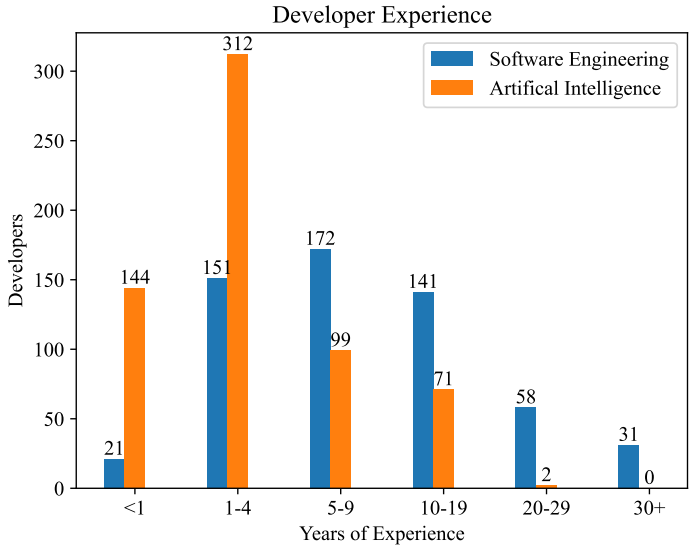
8

Fig. 3. Developer Experience

together and codes that were overbroad could be broken into multiple, more specific codes. This was done when the annotators observed multiple codes frequently being used to describe the same statement in a given answer or when a code was being applied to answers with significantly different meanings, respectively.

Two open-ended questions solicited additional thoughts, comments, or explanations on respondents' answers to a number of different questions. The disparity in responses prevented meaningful analysis with our previous coding method, so we applied a different strategy. One SE researcher read through all 574 responses for both questions and grouped them into categories and sub-categories based on their content by considering the respondent's other answers as context. An additional researcher reviewed and validated these categorization decisions, discussing and resolving disagreements when needed. Lastly, a third author reviewed the quotes extracted for the paper, making sure none were miscategorized, misattributed, or taken out of context. For the closed-ended questions, we aggregated results using descriptive statistics and discussed them to resolve any disagreements. One author manually coded each response for questions with an "Other" option (such as those asking about languages or tools used). Sections 4, 5, and 6 present the results of our analysis.

Where portions of open-ended questions are quoted, we corrected grammar and spelling errors in participant responses for readability. Deletions and additions to text were occasionally made for clarity or space-related reasons and are indicated by ellipses or brackets. Responses have been attributed to survey participants using anonymous IDs for traceability (*e.g.*, $R_{90}$).

### 3.4 Participant demographics and background

Respondents self-selected the nature of their development work from a list that included proprietary software development (331), open source (247), and academia (117). The total number of responses to this question is greater than the number of complete and valid responses (574) because respondents could select more than one option. Respondents came from six continents, as indicated in Figure 4, and 73 different countries, the top ten of which are shown in Table 3.

Table 2. Programming language use reported by developers

| Top-20 Programming Languages | | | |
|---|---|---|---|
| Python* | 407 | Ruby | 11 |
| JavaScript* | 327 | Bash | 10 |
| C / C++* | 161 | Kotlin | 9 |
| Java* | 118 | Swift | 9 |
| Golang* | 106 | Dart | 6 |
| PHP* | 84 | Julia | 5 |
| C#* | 77 | Shell | 4 |
| Rust | 55 | Scala | 3 |
| R* | 41 | Elixir | 3 |
| TypeScript | 25 | Lua | 3 |

*Language was listed as a multi-select option



Fig. 4. Continent of Origin*

| Country | Count | Perc. |
|---|---|---|
| United States | 176 | 30.66% |
| China | 49 | 8.54% |
| India | 28 | 4.80% |
| Germany | 26 | 4.53% |
| United Kingdom | 26 | 4.53% |
| France | 24 | 4.18% |
| Brazil | 20 | 3.48% |
| Canada | 20 | 3.48% |
| Italy | 14 | 2.44% |
| South Korea | 13 | 2.26% |

Table 3. Top-10 country breakdown*

*Some respondents reported working in more than one country/continent

Of the 574 complete and valid responses we obtained, 554 developers (96.5%) indicated that they were generally familiar with using GenAI for producing source code. Of those, 497 (89.7%) had incorporated such tools into their development workflow, as indicated in Table 1. Respondents also used a wide array of programming languages, 48 in total. The most popular languages used were Python, JavaScript, C/C++, and Java, as indicated in Table 2. Respondents had experience with many different AI tools, including GitHub Copilot (73.6%), ChatGPT (84.9%), Bard/Gemini (19.1%), and Claude (10.3%), as well as several open-source models (7.8%); more than 60 tools and models were reported in total, as indicated in Table 4.

### 3.5 Follow-up Interviews

To supplement the survey results, we conducted follow-up interviews. These interviews were designed to establish additional context, obtain answers to questions that arose during the analysis and aggregation of survey responses, and perform a more in-depth exploration of developers' perceptions.

All survey participants were asked at the conclusion of the survey if they would be willing to be contacted for a follow-up interview. Those who consented by supplying their email address (382, 66.5%) made up our potential interview participant pool. Two authors then selected 22 potential interview candidates from this pool based on various criteria, including experience, geographic

Trevor Stalnaker, Nathan Wintersgill, Oscar Chaparro, Laura A. Heymann, Massimiliano Di Penta, Daniel M German, and Denys Poshyvanyk

Table 4. GenAI tools reported by developers

| Top-20 Reported AI Tools | | | |
|---|---|---|---|
| GitHub Copilot* | 366 | Phind | 10 |
| ChatGPT-4* | 304 | Tabnine | 10 |
| ChatGPT-3.5* | 268 | Mistral | 9 |
| Bard* | 95 | Ollama | 7 |
| Claude (Anthropic) (any version)* | 51 | JetBrains AI | 7 |
| Amazon CodeWhisperer* | 33 | Llama / Llama 2 | 7 |
| Open-source models generally | 19 | DeepSeek Coder | 6 |
| Codeium | 16 | Mixtral | 6 |
| Perplexity AI | 12 | In-house or custom model | 6 |
| Codellama | 11 | Codium | 5 |

*Tool/model was included in a list of choices; all others were provided
by respondents to elaborate on their choice of "Other"

jurisdiction, and indicated views/perceptions on key legal issues, with the goal of interviewing participants representing a diversity of views. The full list of interview selection criteria can be found in our replication package [4]. Of the 22 candidates, seven responded to schedule an interview. The seven interviewees represented four continents (North America (3), Europe (2), South America (1), and Africa (1)) and a variety of backgrounds in software development, including working as a developer on proprietary software (6), contributing to OSS (4), and working in academia (1).

The duration of the interviews was between 20 and 30 minutes. The interviews were conducted using the Zoom video conferencing platform. Interview sessions were recorded and transcribed using OpenAI's Whisper large-v3 model [107, 123] to facilitate conversation analysis. A set of shared interview questions was iteratively derived through the collaboration of six SE researchers and one law researcher. This list of questions is available in our replication package [4]. Researchers were also free to ask follow-up or clarification questions during the interviews.

Each interview was conducted jointly by two SE researchers, and the recording and transcript for each interview was independently reviewed by one interviewer to ensure the accuracy of the transcript and extract the relevant portions of each response, assigning topic labels to each one. These quotes and topic labels were hosted in a shared document accessible to both researchers, allowing the reuse of topic labels. Each researcher analyzed roughly half of the interviews. Afterwards, the two researchers jointly reviewed the labels assigned to each response to verify the accuracy and completeness of the analysis, discussing and reaching a consensus on the framing and application of topic labels. A third researcher reviewed the quotes extracted for the paper, making sure none were miscategorized, misattributed, or taken out of context. For simplicity and clarity, when we refer to quotes from follow-up interviews, we use the notation $RI$ and describe participants by their assigned survey IDs.

Interview participants $RI_{31}$, $RI_{118}$, and $RI_{190}$ have experience in proprietary software development. $RI_{118}$ specifically works as a cyber security manager. $RI_{516}$ has contributed to open-source projects. $RI_{39}$, $RI_{149}$, and $RI_{431}$ work as closed-source software developers but have also contributed to open-source projects. $RI_{431}$ also has some experience writing software in an academic environment.

## 4 RQ$_1$ USE OF GENERATIVE AI TECHNOLOGIES FOR SOFTWARE DEVELOPMENT

**RQ$_1$** aimed to discover how developers are using GenAI technologies for developing software. Although the developers surveyed had a range of perceptions about AI and copyright issues, most seemed to share the belief that the use of GenAI is now a regular part of everyday development

activities. As respondent $R_{90}$ commented, "[D]evs not using AI is like accountants not using Excel in the 80s." In a follow-up interview, $RI_{149}$ commented, "If one of my engineers isn't using a code generation model, I'm going to fire him pretty soon. Because there's very little virtue in doing useless grunt work. [...] I actually expect it to be used on every part of the project. If it isn't, that's the exception. In fact, I want to know why [it wasn't used]." In what follows, we present and discuss the findings for this RQ. A summary of key findings can be found in Table 5.

## 4.1 Uses, benefits, and challenges of GenAI

Before we can properly understand developers' thoughts and perspectives on copyright and other legal issues pertaining to GenAI, we must first explore the context in which they use these tools. In this section, we consider the usage scenarios identified by developers, the benefits they have derived from this use, and the challenges that they face.

*4.1.1 GenAI usage in software development.* Almost all respondents (554) stated that they were generally familiar with the use of GenAI to produce source code. Of these, the vast majority—497 (89.7%)—indicated that they had in some capacity incorporated such tools into their development workflows. The remaining 57 (10.3%) offered several reasons for not regularly using GenAI, including legal risk (8), personal preference (8), the lack of significant use cases warranting its use (6), reluctance to send sensitive information to a remote server (6), constraints imposed by an employer (6), and limited trust in the tooling (5).

Among those who had a personal preference against using GenAI, some developers conveyed the sense that GenAI encourages monotony or lack of effort. This was somewhat surprising because, ideally, GenAI is used to automate repetitive and boring tasks, as previous work has indicated [112]. $R_{188}$ stated "I like to think for myself. [W]ithout [that], coding becomes boring and a lot of reading instead of reading and writing." $R_{339}$ put it this way: "It's better to write code on my own than to review code generated by AI." $R_{276}$ reported similarly, "I've tested incorporating GitHub's [C]opilot into my workflow and caught myself many times waiting for the autocompletion to kick in so I could just press 'Tab' and go from there. Soon I realized I was getting lazy and not thinking much, which I didn't enjoy." The full list of reasons put forward by developers can be found in our replication package [4].

> **Finding 1:** While the vast majority of developers surveyed (89.7%) were using GenAI tools, there were many reasons why developers might not incorporate GenAI into their development workflows, including perceived legal issues, personal preference, information security concerns, and limitations placed by employers.

Of those familiar with GenAI tools, relatively few developers (32) said that their organization did not permit the use of GenAI for coding. These developers reported their understanding of the restriction. Twelve (37.5%) stated that such usage could compromise sensitive or proprietary information. Ten (31.3%) indicated that their organizations were concerned about the legal implications of using such tools. Four (12.5%) suggested that their employers were slow to adopt new technology and did not fully understand its benefits. Three respondents (9.4%) pointed to the fear of security issues, such as incorporating insecure code or designs. Notably, these responses could reflect overlapping concerns: compromising proprietary information, for example, could also raise legal concerns.

> **Finding 2:** Of those respondents who reported working at organizations that disallowed the use of GenAI, many indicated that they believed the restriction was based on fear of proprietary information leakage and the potential legal challenges such usage might bring. Some developers also perceived that their employers were slow to adopt new technology.
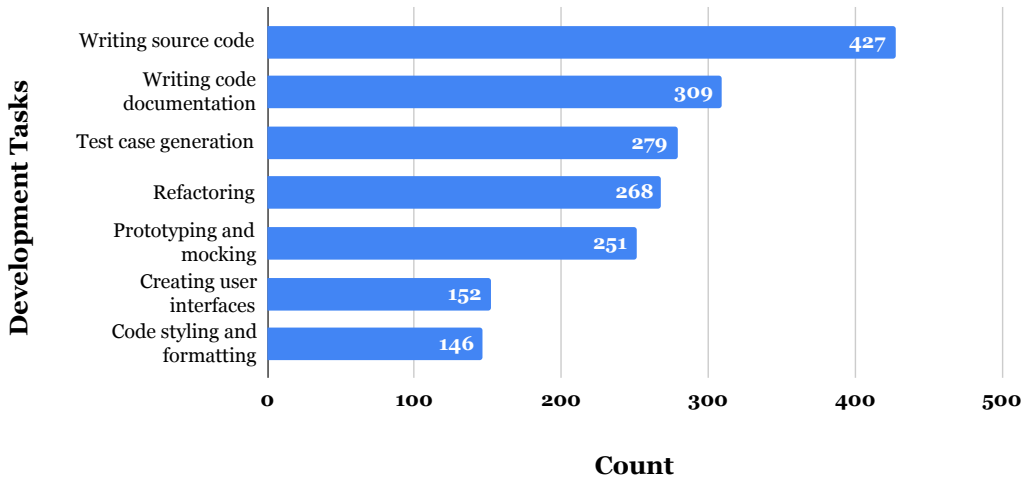
Fig. 5. Software development tasks where developers reported using GenAI

Respondents have used GenAI for a variety of software development tasks (as seen in Figure 5), including writing source code (427), writing code documentation (309), test case generation (279), refactoring (268), prototyping (251), creating UIs (152), and code styling/formatting (146). Developers also mentioned that they used such tools for debugging (12), reference or learning (12), code explanation (8), troubleshooting (6), repository management (5), and brainstorming (4).

*4.1.2 GenAI tools used by developers.* A wide variety of GenAI tools, including open-source, open-weight, and proprietary models, were used by participants. In total, developers reported using 64 distinct, named GenAI tools, the top 20 of which can be seen in Table 4. Developers were encouraged to identify all the tools that they had used in their work. Of these, the top six (which were also provided in a list of choices, along with an "Other" option) were: GitHub Copilot (366), ChatGPT-4 (304), ChatGPT-3.5 (268), Bard (95), Claude (51), and Amazon CodeWhisperer (33). The next largest group (19) reported that they used open-source/open-weight models, with several such tools/models listed by multiple respondents, including Codellama (11), Mistral (9), Ollama (7), Llama2 (7), and Mixtral (6). Developers also used other proprietary models and services, including Codeium (16), Perplexity (12), Phind (10), and JetBrains AI (7), and six respondents indicated that they used custom-trained or in-house models for their development work. A full list of the tools developers reported using can be found in our replication package [4].

Of the tools listed by respondents, 20 were designed to be integrated into an IDE. Many of these tools provide autocomplete suggestions for developers based on the context of surrounding code, but some also include Q/A chat functionality. Developers also reported using 18 distinct web-based tools. These tools, such as ChatGPT [100], Perplexity [104], and Hugging Face Chat [62], offer great convience, but require the user to supply more context. Respondents also indicated using four different model infrastructures: Ollama [98], Oobabooga [99], llama.cpp [82], and Mozilla llamafile [59]. These infrastructures provide a streamlined means to run and train open-weight models locally. In total, 15 distinct open-weight models/model families were also identified by respondents. (We note that the prevalence of open-weight models and open-source tools in this space has likely increased since our initial survey with the subsequent release of models such as Llama3 [90] and the open-sourcing of XAI's Grok [126].) Some developers also noted the use of
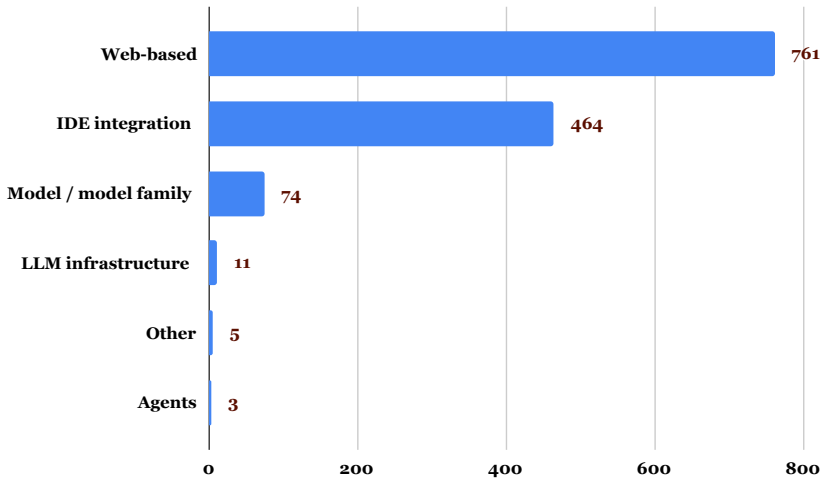
Fig. 6. Instances of tool category uses reported by developers

GenAI agents like gpt-engineer [50] and metagpt [58], which allow models to self-prompt with the goal of accomplishing more complex tasks. Lastly, four GenAI tools which were not strictly for software development, such as Figma [39] and GrammarlyGo [51], were mentioned.

Figure 6 shows the reported tool usages across the categories previously discussed. Since developers selected all the tools they had used for development tasks and many of those tools might fall into the same categories, we can't generalize to reach a conclusion about any developer preferences. We can also say little about the frequencies with which each tool type is used. That said, we see that there were 761 reported usages of web-based tools. There were only 464 reported usages of IDE-based tools, despite them constituting the largest number of distinct tools. This might suggest that developers were more likely to experiment with various different web-based tooling solutions, especially during the early days of GenAI for software development. There were 74 reported usages of open-weight or locally hosted models, but only 11 instances of reported corresponding model infrastructure. It could be that many developers using open-weight models neglected to mention the tooling they used to run them or that they relied on web-hosted solutions such as Hugging Face Chat.

> **Finding 3:** The ecosystem of GenAI tools used by developers was varied and extensive, comprising both open-source/open-weight and proprietary solutions. Developers primarily made use of web-based tools and those that could be integrated into their IDEs.

*4.1.3 Benefits derived from GenAI usage.* Respondents who had used GenAI as part of their workflow identified several benefits from its use. The greatest benefit mentioned by 288 developers (58%) was increased productivity, faster development, and efficiency. Other benefits included not having to spend time writing simple or boilerplate code (88), quick prototyping (42), debugging (37), faster, more accurate documentation (29), summarizing existing documentation as an alternative to reading it (28), providing inspiration (16), refactoring (14), use in writing test cases (5), domain logic support (4), and generating test data (4). A full list of the benefits identified by respondents can be found in our replication package [4].

Thirty respondents identified improvements in code quality and optimization over what human developers could achieve. $RI_{118}$ noted, "When you couple the lack of training [of the average

developer] with generally one to three years' experience, and the fact is that GPT4 or many of these tools are actually better developers than they are. [...] I would almost feel better if [a code comment] said I wrote this and then I got a code review from GPT4." $R_{74}$ stated that GenAI's ability to "notic[e] optimizations or cleaner ways of writing something like a custom function is quite useful." Some respondents, like $R_{466}$, noted that the "[generated] code [is] more reusable and better," and $R_{669}$ noted that the output of code generation models was "following best [coding] practices."

Similarly, developers related that GenAI could help them write software in unfamiliar languages (28) and libraries (27). During a follow-up interview, $RI_{190}$ described how they would "basically just writ[e] the code in Typescript and then just tell [ChatGPT] convert to Liquid" and that "otherwise [they] would probably quit because Liquid is incredibly painful." Using GenAI specifically in the capacity of porting or translating code was mentioned by six respondents. $R_{286}$ found that GenAI "was helpful to migrate code from a middleware platform to Javascript based development." $R_{90}$ summarized this benefit: "The syntax of various languages is [why it] can take hundreds of hours of use to establish [familiarity] in a new language, but AI significantly reduces the ramp-up time."

Along the same lines, several respondents noted that GenAI helped them to be better developers by enhancing their understanding of unfamiliar code. For example, nineteen respondents stated that GenAI provided "assistance in [the] understanding of code snippets" ($R_{47}$). $R_{404}$ also explained that "[GenAI] is also very good at providing semi-formal abstract explanations of the topics at hand and it provides consistent analogies which makes things very smooth to follow and comprehend." But GenAI models are not limited to explaining source code or concepts. $R_{238}$ noted that they could also be used for "commit message generation [and ex]planation of changes in commit (summarization)."

> **Finding 4:** The benefits developers reported from using GenAI tools for coding go beyond increasing productivity and avoiding repetitive, boring tasks. They also included understanding code better and improved code quality and code optimization, which some human developers might struggle to achieve.

*4.1.4 Challenges and shortcomings encountered in using GenAI.* Respondents also identified several challenges arising from the use of GenAI, the most significant of which was the generation of unhelpful, unwanted, or broken code (128). Additional challenges included hallucinations (64) (such as recommending non-existent libraries or method calls), problems arising from outdated training data (53), issues arising from tools not having access to the project's wider context (52), which could be caused by insufficient context windows (40), and an inability to use GenAI to solve complex problems (43). $R_{291}$ noted, "[A]s issues get more complex, code generation has many gaps." $R_{148}$ elaborates: "AI can do simple code pretty well, it's not that capable [of generating] more complex code, it's incapable of generating well-architected code outside of extremely simple patterns." $R_{33}$ provided an example: "In my experience, a simple Python/Django project that required some Form customizations was just not possible for any AI tool at my disposal. The logical part, the thinking, and the experience gained from years of coding real-world applications and systems design are [not comparable to] an AI code generation tool like GPT or Gemini." The full list of challenges reported by developers is included in our replication package [4].

Twenty-five respondents also noted that while the ability to write good prompts is essential to effectively using GenAI tools, it can present a significant challenge. $R_{303}$ reported that "AI cannot complete the majority of tasks by itself, even simple ones[,] without time spent [by the developer on] prompt engineering." As $R_{378}$ put it, "Models need 'hand-holding' instructions [in order] to generate code and debug it." But when crafting these prompts, "[q]uestions must be asked precisely," as $R_{135}$ commented. $R_{432}$ elaborated, "Vague prompts lead to vague answers, so you have to be

specific with what you want." $R_{72}$ summarized all these sentiments: "The outputs [of GenAI tools] are only as good as the inputs."

Writing detailed and specific prompts isn't always an easy task, however. $R_{18}$ said, "It is hard to communicate what I want or what I'm trying to describe to the model." $R_{83}$ echoed this concern, stating, "Sometimes I am unable to describe what it is that [I] want[,] only to find out [that] there existed vocabulary that I should have used[, but] [t]his issue seems less frequent recently compared to when [I] first started using [GenAI] for coding." If done correctly, however, prompting was also seen as a way to overcome other limitations of the technology. $R_{240}$ said, "There is a limit to the complexity that can be generated stand[ing] alone, but a human crafting intelligent prompts can still architect complex solutions using AI."

> **Finding 5:** Aside from obvious problems such as possible hallucinations, GenAI for coding still suffers from challenges related to lack of suitable contextual information available to the model at inference time, the inability for such models to solve complex, non-trivial tasks such as highly-specific functionality which conforms to a given architecture, and the need to provide good, well-designed prompts.

## 4.2 Current practices for documentation and legal compliance when using an AI tool

*4.2.1 AI-generated code review and compliance.* Given the potential challenges associated with the use of GenAI, mitigation strategies relating to vulnerability and compliance analysis become increasingly important. In our follow-up interviews, we asked participants additional questions regarding their workflows. (Section 3.5 provides demographic information on participants, including development experience background.)

As $RI_{118}$ noted, when code is not generated by the developer, review can be a challenge: "[T]hat's another big part of it is just going through line by line, but it's like a more passive mindset, which I find a bit challenging compared to actually writing the code, so it can be harder to actually spot what's wrong with it. So the tests become really important." This is particularly true with respect to edge cases, "because the tools are actually quite good, but unimaginative and have no domain, [they] have no context."

Some developers reported an overall review but not focusing on vulnerabilities. $RI_{190}$ said, "[I don't do] any checks for licensing or vulnerabilities, but usually I read the code before I commit it and it's also reviewed in a pull request." $RI_{39}$ offered a similar approach: "I think it's mostly that I read the code and try and evaluate first of all how dangerous it is to run."

Others believed that AI itself could assist with reviewing AI-generated code. $RI_{149}$ told us that they "ask the GenAI itself to review the code, perhaps using a different model [...] [and] ask it to give me gotchas or things to look out for." $RI_{31}$ tentatively agreed but acknowledged the challenges: "I've thought about this [...,] and [if] it's something where the answer is using machine learning to check whether or not the code that you've used appears, and would put you at legal risk, [...] you're relying on a large language model to check whether your large language model puts you at legal risk."

Some developers, however, employed manual checks, particularly for license compliance. $RI_{516}$ described a process that involved "search[ing] on GitHub because that's the largest repository, but [...] also check[ing] inputs in search engines, finding other places like GitLab or CodeBerg or Bitbucket or other repositories. And if the code is in a visible repository, it should show up." $RI_{431}$ would "try to do my own research on websites such as Stack Overflow or GitHub." $RI_{431}$ emphasized the manual nature of their compliance review: "I don't really have any tool [to verify provenance], so that's why I'm doing this step of going to a website such as Stack Overflow and trying to see if the code which has been given to me by the model has already been proposed somewhere else on the Internet." But $RI_{516}$ noted that these manual processes are limited, because "if it's not open
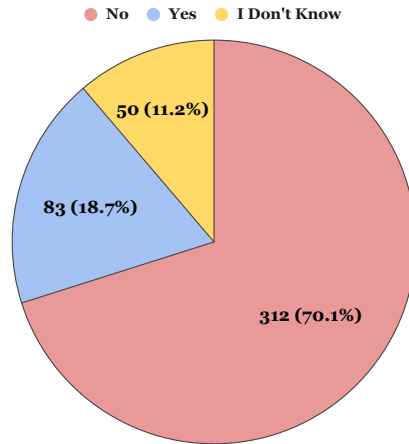
Fig. 7. Developer awareness of documentation processes for GenAI-related tasks within their organizations

source, I probably have no way to or means to check it because closed source is not generally available [...] in a search or anything like that."

Some developers stated that they avoided compliance checks on AI-generated code for fear of opening themselves or their organizations to liability. $RI_{39}$ explained, "I have gotten advice from IP lawyers at a previous job [...] to not look at patents for anything that I'm working on. Because if there is evidence that I have seen a patent that we then infringe on, then it's willful infringement, which can mean triple damages or something like that. [...] I could imagine a case where it might make sense to not track down which lines of code were generated by a LLM, because then you might have less of a willful infringement case or something like that."

> **Finding 6:** Among respondents, GenAI code review is typically done manually or by using other GenAI to check for code correctness and compliance. Some developers are concerned that the process of conducting a review may itself expose them or their employer to liability.

*4.2.2   Documentation of GenAI usage.* One might expect that organizations that permit the use of GenAI for code development would ensure that developers were aware of processes to document its use. However, of the 445 developers indicating that their organization (partially) permitted the use of GenAI, only 83 (18.7%) indicated that they were aware of such a process, and of the remaining group, 312 (70.1%) said that there was no process and another 50 (11.2%) were unsure if a process existed. $R_{325}$ offered one view of why a process was unnecessary: "I did not use anything that cannot be googled anyway—but we never thought about having to document or question code we found through Google." Whether this means that many organizations have no process for documenting the use of GenAI or that at least some developers incorrectly believe that no documentation process exists at their organization, the result may be that developers do not have a full understanding of the possible legal risks they or their organizations face.

> **Finding 7:** Most developers in our study (70.1%) were unaware of any process for documenting GenAI usage within their organization. To the extent that organizations have processes in place for documenting the use of GenAI, they may not be sufficiently educating their employees about those processes.

We further asked the 83 respondents who reported that their companies had processes to document their AI usage about what information was collected during the process. The majority of
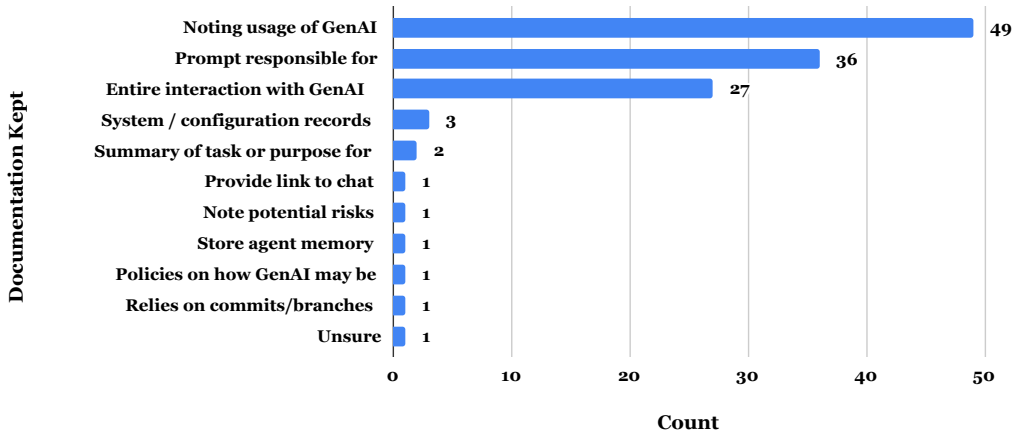
Fig. 8. Information collected / stored during documentation processes

responses (49, 59.0%) indicated that documentation consisted of noting the use of an AI tool, with 36 respondents (43.4%) indicating that the prompt was also documented. In 27 cases (32.5%), the entire interaction with the code generation tool was documented. Other documentation strategies included keeping system/configuration records (3), summarizing the purpose of using AI tooling (2), providing links to chats (1), documenting known risks associated with the use of AI (1), storing the AI agent's memory (1), and relying on version control (1).

> **Finding 8:** The most common documentation process among respondents was to note the fact of GenAI usage and/or to document the prompt used to acquire the generation. However, nearly a third of respondents (32.5%) indicated documenting the whole interaction, potentially leading to reproducibility.

Some respondents indicated that one motivating factor behind documentation, at least for developers, was to explain the cause of a potential failure rather than to provide provenance. $R_{190}$ noted, "I usually write a comment 'written by [C]hatGPT' or 'partly written by [C]hatGPT' to denote that if it breaks, I would maybe have written it better but used [C]hatGPT to speed up the process." $RI_{431}$ echoed this sentiment during a follow-up interview: "But when I do use artificial intelligence to generate code [...], I will try to notify in a code comment or in a comment inside the documentation that this part needs to be rewritten or at least tested [...] [so that we can] change anything that wouldn't work in our current case." $RI_{118}$ suggested that the typical training at their place of employment does not anticipate this type of activity: "[Developers are] allowed to start using the tools and there's guidance on where and when they can't, but not in terms of attribution or [adding] warning[s that say,] 'Hey, here be code generated monsters.'"

*4.2.3 Review of the Terms of Service for AI Tools.* The Terms of Service (ToS) for GenAI tools lay out how developers may interact with models by describing permissible and impermissible behaviors. These ToS may contain language that prevents the use of the service for any illegal activity, forbids the use of model output as training data for another, competing model, or, for open-weight models, provide stipulations on what is allowed when fine-tuning. The ToS may also include stipulations that data provided as prompts by a user may be used as training data for future models. In short, it

Trevor Stalnaker, Nathan Wintersgill, Oscar Chaparro, Laura A. Heymann, Massimiliano Di Penta, Daniel M German, and Denys Poshyvanyk
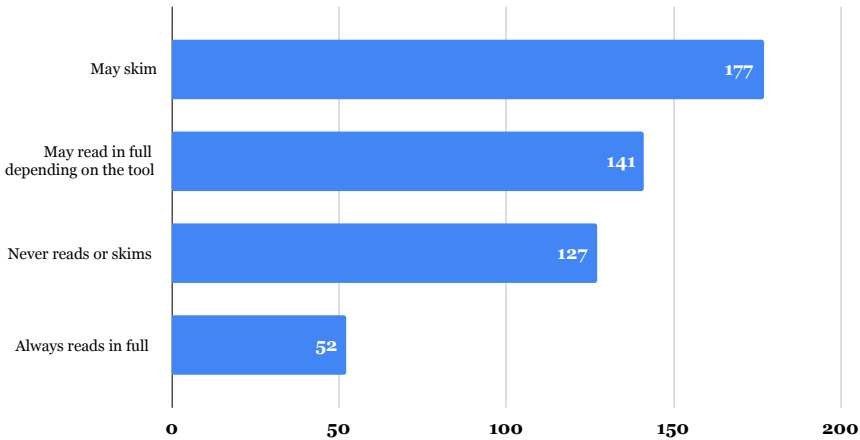
18

Fig. 9. Developer approaches to Terms of Service (ToS) for AI Tools

is important for both users and model fine-tuners to understand the terms that they are agreeing to before using a model.

Developers' views on the ToS of GenAI tools varied widely. A little over a third of respondents who said that they incorporated AI into their workflow (177, 35.6%) indicated that they never read terms in full but may skim them looking for relevant or problematic provisions. A similar percentage of respondents (141, 28.4%) stated that they sometimes read the ToS in full, depending on the tool they are using. An additional 127 respondents (25.6%) claimed that they never read the ToS for the tools they use. The smallest group, 52 participants (10.5%), told us that they always read the ToS in full.

In follow-up interviews, respondents who tended not to read the ToS provided their rationale. For one participant, the decision was based on other priorities. $RI_{190}$ explained, "[As to] why I don't read it, it's because I don't have time. We have features to ship, and I work at a startup, so we need to find this Holy Grail, which is product market fit. [...] It doesn't feel like reading the Terms of Service will help us, so I don't do it." $RI_{431}$ indicated that the main reason they do not read the ToS is "quite often the length of the text."

Still other respondents indicated that they relied on others to read the documents and alert them to relevant provisions. Respondents who worked within organizations stated that they trusted that their organization's legal department would take responsibilty for compliance. $RI_{39}$ stated, "I personally don't worry about [licensing issues related to generated code] because the legal department where I work has already authorized [the] use of [the LLMs], and it's kind of their problem from my perspective."

Others noted that they relied on discussions on blogs and social media to identify problematic provisions. $RI_{190}$ stated, "If it would blow up on the internet that the terms are bad, then I would kind of hear about it." $RI_{31}$ indicated they have typically "already heard about pros and cons" of using certain models and described one rule they were "trying to stick by": "I know that some models and services just copy and paste what's already been created [...] It's stealing work from someone. I don't really agree with that philosophy." $RI_{31}$ also used models' ToS and social media to identify models that purportedly trained on unauthorized data: "You can review their Terms of Service, and look at their blogs and so on. The other piece is just honestly checking Reddit [and] seeing what other people are talking about." $RI_{516}$ referred to the website Terms of Service; Didn't Read (https://tosdr.org) which "[does] a summary of the main points of each Terms of Service and

privacy policy on the services that have already been analyzed. It's analyzed by a community [...]. The main points of concern [...] are marked in red, so you need to watch out for them."

The group of respondents who stated that they did read the ToS for tools they used reported focusing on how the tool developer acquires and uses both personal data and training data, at least for some tools. $RI_{516}$ explained, "I always take a look to see what [GenAI tools] do with the data, the risks involved in inadvertently disclosing personal or private information into the data." $RI_{31}$ expressed similar concerns, "The training set, making sure that the training set is—they haven't trained it on stolen code. And then the other one is making sure that I can run the whole thing locally and maintain privacy." Some others, for example, $RI_{190}$, took a more pragmatic approach, stating, "[The major red flags] would have to be something that could bring me or the company I work at into big trouble, and that would [have to] be more trouble than that we probably won't exist in a year because our product isn't good enough, so it would have to be pretty big."

> **Finding 9:** Even though some provisions in an AI tool's Terms of Service may be important to developers and their work, developers do not always thoroughly read the Terms of Service for the tools they use, choosing instead to skim the document, assume that their organization has read the document, or rely on information and validation from others in the development community.

*4.2.4 Copyright/legal training.* Of the 574 respondents, only 68 (11.9%) indicated that they had any formal training in copyright law and/or the legal implications of using code generation models. Of these, 28 (41.2%) sought the training out on their own. Sixteen respondents (23.5%) indicated that the training was provided by their employer, and 24 (35.3%) described their training as some combination of self-initiated and employer-provided. (Some respondents interpreted the question to refer to GenAI training more generally rather than training that focused on legal implications.)

In our follow-up interviews, respondents provided further information about the nature of the training reported by developers. $RI_{431}$ described their training as "two courses [during my master's degree about] [...] everything related to the law about anything digital." By contrast, $RI_{31}$ experienced a dearth of resources: "I started looking for training not just for myself but for our team, and it wasn't something where I could really find a whole lot. [...] There aren't a lot of resources out there for that because everybody's still figuring out the legal implications."

> **Finding 10:** Relatively few respondents (11.9%) reported having undergone any formal training on the legal implications of GenAI. Some struggled to find suitable available resources.

## 5 RQ$_2$: DEVELOPERS' PERCEPTIONS OF LICENSING AND COPYRIGHT ISSUES EMERGING FROM THE USE OF GENERATIVE AI

**RQ$_2$** sought to determine developer's beliefs and understandings regarding the licensing and copyright issues that arise from the use of GenAI tools. Here we explore developers' thoughts on current litigation and the prospect of regulation, sentiments regarding the use of their own code in model training data, opinions on the ownership of GenAI output and prompts, and lastly, the perceived risk of copyright infringment. A summary of key findings can be found in Table 6.

### 5.1 Perceptions about litigation and regulation

We asked respondents whether they believed that litigation would result in limitation on the use of GenAI for software development. (Figures 10 and 11 show respondents' answers to nine five-point Likert-scale [102] questions.) Of the 554 respondents who indicated a familiarity with AI tools, 205 (37.0%) somewhat agreed with the view that current litigation will result in limitations placed on the use of GenAI models for software development, with 76 (13.7%) respondents strongly agreeing.

Table 5. Main findings for RQ1: Use of GenAI for software development

| | **4.1 Uses, benefits, and challenges of GenAI** |
|---|---|
| | 4.1.1 GenAI usage in Software Development |
| F1 | While the vast majority of developers surveyed (89.7%) were using GenAI tools, there were many reasons why developers might not incorporate GenAI into their development workflows, including perceived legal issues, personal preference, information security concerns, and limitations placed by their employers. |
| F2 | Of those respondents who reported working at organizations that disallowed the use of GenAI, many indicated that they believed the restriction was based on fear of proprietary information leakage and the potential legal challenges such usage might bring. Some developers also perceived that their employers were slow to adopt new technology. |
| | 4.1.2 GenAI tools used by developers |
| F3 | The ecosystem of GenAI tools used by developers was varied and extensive, comprising both open-source/open-weight and proprietary solutions. Developers primarily made use of web-based tools and those that could be integrated into their IDEs. |
| | 4.1.3 Benefits derived from GenAI usage |
| F4 | The benefits developers reported from using GenAI tools for coding go beyond increasing productivity and avoiding repetitive, boring tasks. They also included understanding code better and improved code quality and code optimization, which some human developers might struggle to achieve. |
| | 4.1.4 Challenges and shortcomings encountered in using GenAI |
| F5 | Aside from obvious problems such as possible hallucinations, GenAI for coding still suffers from challenges related to lack of suitable contextual information available to the model at inference time, the inability for such models to solve complex, non-trivial tasks such as highly-specific functionality which conforms to a given architecture, and the need to provide good, well-designed prompts. |
| | **4.2 Current practices for documentation and compliance** |
| | 4.2.1 AI-generated code review and compliance |
| F6 | Among respondents, GenAI code review is typically done manually or by using other GenAI to check for code correctness and compliance. Some developers are concerned that the process of conducting a review may expose them or their employers to liability. |
| | 4.2.2 Documentation of GenAI usage |
| F7 | Most developers in our study (70.1%) were unaware of any process for documenting GenAI usage within their organization. To the extent that organizations have processes in place for documenting the use of GenAI, they may not be sufficiently educating their employees about those processes. |
| F8 | The most common documentation process among respondents was to note the fact of GenAI usage and/or to document the prompt used to acquire the generation. However, nearly a third of respondents (32.5%) indicated documenting the whole interaction, potentially leading to reproducibility. |
| | 4.2.3 Review of Terms of Service |
| F9 | Even though some provisions in a tool's Terms of Service may be important to developers and their work, developers do not always thoroughly read the Terms of Service for the tools they use, choosing instead to skim the document, assume that their organization has read the document, or rely on information and validation from others in the development community. |
| | 4.2.4 Copyright/legal training |
| F10 | Relatively few respondents (11.9%) reported having undergone any formal training on the legal implications of GenAI. Some struggled to find suitable available resources. |

An additional 109 (19.7%) respondents were neutral (perhaps reflective of the legal uncertainties at play), and the third largest group of 104 developers (18.8%) somewhat disagreed. The smallest group, 60 developers (10.8%), strongly disagreed.

While more respondents believed that litigation would result in limitations on the use of GenAI, developers who believed that litigation would *not* result in limitations were more vocal about their views in their responses to open-ended follow-up questions. Indeed, many respondents indicated that technology would outpace attempts to regulate it, whether through litigation or otherwise. $R_{98}$ commented, "I don't see any limitations coming due to ongoing litigation—yes, it may cause issues for OpenAI and StabilityAI, but you can't put the genie back in the lamp." $R_{507}$ further elaborated:
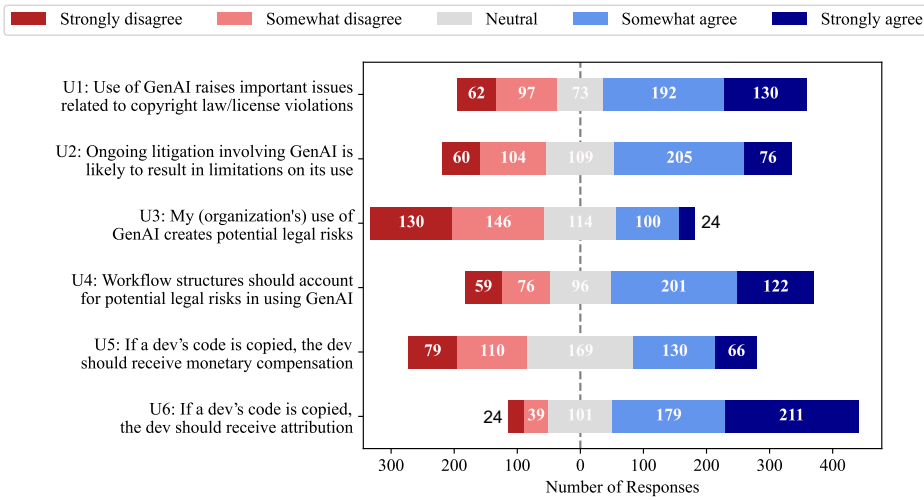
Fig. 10. Developer perceptions on copyright and GenAI

Fig. 11. Developer attitudes towards copyright and GenAI

"Pandora's box has been opened. Even if the current models were legally limited, good luck proving that certain code has been used as training data, and good luck stopping people from running Llama instances locally." $R_{54}$ acknowledged the existence of legal issues but believed that "the momentum that's already been built by OpenAI and all the hype around the tools means there will not be the political capital to get any restrictive regulations in place." $R_{250}$ said, "Generative AI is used by almost all software developers at this point. That ship has sailed. Lawyers should give up on even trying to figure this out."

One respondent ($R_{239}$), after weighing the pros and cons of regulation, concluded that no regulation at all was better than regulation that hampered development: "On the one hand, I know it's sort of a legal gray area right now and ultimately we're going to need legislation. On the other hand, these tools are incredible[,] and I can't imagine going back to not having Copilot. It would be

like riding a bicycle with a flat tire. I would like to see sensible legislation that promotes growth and advancement of these technologies, but I'm worried that we won't get that. I'd rather it continue to be a free-for-all—including, yes, commercial models trained on freely available source code—than have the tools taken away or severely restricted."

$R_{329}$ recommended a cautious approach to regulation and litigation, stating that if "[l]awyers get involved before the technology matures[, we are] putting the carts before the horses. Preemptive copyright protection can only hurt technology innovation." In the same vein, $R_{185}$ said, "I am skeptical about what I consider to be premature efforts to restructure the legal framework around AI until we get a better handle on its implications and potential."

During a follow-up interview, $RI_{118}$ mused on a potential regulatory body: "No one has solved regulatory capture, so [...] I don't know how to solve that one. It would probably have to be like a NIST [National Institute of Standards and Technology] or some kind of body like that. [...] Hopefully [it] isn't run by presidential appointment, so the datasets are not decided by the politics of the day [and not by] an industry rubber stamp body because that doesn't help either. [...] So it would probably have to be some standards organization."

Finally, some respondents conveyed the view that it would not be technology that changed to adapt to the law but rather the law that would change to adapt to technology. $R_{376}$ commented, "The result [...] will probably not be limitations on the technology since open-source is rapidly developing and cannot be genuinely regulated in this fashion. Rather, intellectual property, artistic rights, etc. will need to be recalibrated." $R_{79}$ echoed this sentiment: "The old style copyright will not make any sense as we know it as we are entering the bumpy road to abundance." $R_8$ similarly put the burden on the law to adapt: "The legal side will be resolved by lawyers when they catch up with what these tools actually do." And $R_{182}$ offered an aspirational approach, stating that they "don't think it is wise to approach the question of copyright of AI-generated/assisted code as a matter of whether the generated code is or is not copyrightable based primarily on what has been settled. Rather, we should think about what we want it to be, keeping in mind the benefits copyright is intended to bring to humanity."

> **Finding 11:** While about half of participants anticipated that litigation would result in limitations on the use of GenAI, a vocal minority suggested that imposing such limitations would be difficult if not impossible, and others urged restraint, given the pace of technological developments.

## 5.2 Using one's own code in models' training data or as output

If regulation of GenAI occurs, whether through legislation or litigation, it should ideally take account of the perceptions and practices of stakeholders, including software developers. Our study, therefore, sought to learn respondents' views regarding the use of code as training data for generative models.

*5.2.1 Sentiments regarding the inclusion of one's own code in training data.* Of the 554 respondents familiar with GenAI tools, the largest proportion (40.6%), as seen in Figure 11, indicated that they would be neither pleased nor displeased if the code they developed was included in training data without their permission. However, open-ended responses indicate that the reality is context-dependent. If the work was open-sourced, respondents generally indicated that they would be pleased or even flattered to have their code included in training data as a way in which they could contribute to the greater good. $R_{144}$ expressed this sentiment: "If it's just a class, or a few elements from my code, I would actually be pleased and proud that my code is having [an] impact on other peoples' work and on the advancement of science in general." $R_{67}$ elaborated, saying, "If the code has been used [in training data] with my implicit (*e.g.*, through BSD-licensed code) or

explicit (*e.g.*, through my company's or my permission) then I would be happy to have contributed. Otherwise [my] answer becomes 'displeased.' " (Some respondents, however, expressed concerns about training models on open-source code given the quality of the training data. $R_{196}$ stated, "[M]ost of the code I write is bad or middling, and the lack of intentional design I've seen in the code of open-source projects [...] really worries me [...] [in that] we can just automatically generate unclear and potentially insecure code.")

Conversely, developers indicated disappointment about the idea of their proprietary code being included in training data. $R_{82}$ said, "I'm perfectly OK if a company trains their model on MIT [licensed] projects that I have on GitHub, but GitHub training their model on [proprietary] code is a bit different." Similarly, $R_{98}$ indicated that for private projects, "I wouldn't be happy if the code was used for LLM training mostly because some of these repositories may contain committed secrets," and for work projects, "the code in these is in private repositories and everything that's inside is a corporate secret. If that code is used for LLM training, it could harm my employer by exposing secrets and algorithms we use, which would be extremely displeasing."

Some developers were also concerned that the unauthorized use of work as training data for GenAI would disincentivize human creativity, consistent with the traditional justification for copyright law in the United States. If users increasingly seek solutions from LLMs rather than directly from the original project, developers may be reluctant to continue to contribute to the open-source community, causing the well of open-source training data to eventually run dry. As $R_{148}$ put it, "[This] is more of a case of how to sustain OS[S], than [it is] a case of code being copied." Indeed, $R_{245}$ said, if developers are not incentivized to create code, "we will end up with a gradual decline in the quality of all intellectual property as current AI systems are not creative enough to produce IP that is significantly different from the training material." $R_{514}$ offered a similar view: "The real issue to me is that those engines are [parasiting off] websites like [S]tack [O]verflow [and] actually lowering the amount of money they make. If these websites stop [...] produc[ing] good answers that [AI tools] can feed [off of], they [the AI tools] all will die." $R_{604}$ even speculated that "[c]ode generation tools might encourage creative professionals to hoard knowledge, defeating the original purpose of copyright: enabling protected[,] limited-time monetization followed by release to the public domain."

> **Finding 12:** Developers' views on having their code included in models' training data was context-dependent. If the code was open source, many developers were indifferent or even pleased to have their work included; this was not true for proprietary code. Developers also expressed concerns about the quality of code incorporated into the training data and whether increased use of GenAI tools would disincentivize developers from contributing to the OSS community.

*5.2.2  Conditions around use of one's code in training data.* If using another's code in training data requires permission, the question then becomes what conditions might be imposed, including monetary compensation and attribution.

Some respondents argued that the inclusion of proprietary code in training data warranted monetary compensation. Relatedly, many believed that the developers of open-source software had waived their right to any compensation. For some respondents, the answer depended on whether the owner of the AI model profited from the use of the data. $R_{159}$, for example, stated, "[I]f my code is free, and the AI model is paid, I feel wronged and deserve compensation." $R_{613}$ offered a similar view: "[T]he developers whose code these big corporations use to train their models must receive some amount of compensation in terms of free credits, [GPU] time[,] etc. on their platform."

$R_{433}$ took a broader view: "If monetary compensation were required to train these models, it would put smaller AI development companies out of business and restrict the currently thriving sector to only the Tech Giants. No one else would have the [means] to not only acquire the data but also to do the training."

> **Finding 13:** Some respondents suggested that training data contributors should receive monetary compensation or other benefits, such as free access to the trained model or free credits for the platform. However, forcing AI companies to provide monetary compensation to all contributors of the training data, while seemingly fair, would place an additional hurdle potentially blocking smaller AI companies from entering the market.

For many respondents, it was also important for licensing terms to be respected, including attribution for their contributions or (for restrictive licenses) reuse under the same terms. Modes of attribution, as respondents suggested, could be as fine-grained as a citation provided along with the generated code or as coarse as a manifest that included links to all the training data. In the view of respondents, embedding attribution into the AI model would serve two functions: (1) respecting the rights of developers and (2) helping users and organizations ensure that they are complying with the licenses of code they use by providing provenance of the output. As to the latter, $R_{456}$ imagined a situation where an organization using GenAI unintentionally reproduced GPL-licensed code verbatim, introducing license compliance issues into the project. At least two respondents felt that responsibility for avoiding such problems rested with the models. $R_{378}$ said, "The onus is on the LLM creators to make sure [the output] includes the source" because "[t]he users are not at fault[...]." $R_{665}$ noted, similarly, "Developers should not bear the burden for [...] copyright [...] concerns. Tools should not be producing code which could get a developer in trouble."

Some respondents would favor the creation of new OSS licenses that provide specific conditions for inclusion in training data. Such licenses would enable creators to control how and when their works are used for GenAI. As $R_{587}$ put it: "I think that open source developers that hold strong opinions about the use of their code with AI models need to update their licenses to enforce those desires."

> **Finding 14:** Respondents generally wanted licensing terms to be complied with by GenAI models, especially when producing code identical to or very similar to work in the training set.

*5.2.3  Fair use and replicated code.* While some developers, as described above, believed that the use of code as training data requires a license, others considered such use to be a fair use, even if it was not always clear whether respondents were referring to this concept in its legal sense or in a more general sense. For example, $R_{173}$ said, "I believe that AI training utilizing any published content is a form of fair use and should be encouraged rather than litigated. The usage is entirely transformative rather than derivative." $R_{695}$ said, "I think training on copyrighted materials is legal, but their outputs can violate [copyright law] surprisingly easily (even worse when considering trademarks). I've had them reproduce IDENTICAL text to copyrighted sources, which is clearly infringement. AI companies want to not have liability for such infringement, putting it on the users, but that's stupid." For some, the extent of the copying was important, perhaps reflecting U.S. copyright law's concept of *de minimis* copying. $R_{144}$ put it this way: "I feel like if significant pieces of code are copied, only then does it really become an issue for me: I would not be comfortable if someone took multiple files from my code and used them in a proprietary project." $R_{143}$ said, "If a generative AI happens to have copied large swathes of code, yes [the developers of that code deserve attribution]."

Some respondents who believed that the use of content as training data was fair use compared the practice to the process through which humans produce new works based on the knowledge they have accumulated over time. These respondents therefore believed that it should be assumed that code made publicly available is free to use for training, and the burden should be on developers to indicate otherwise. For example, $R_{26}$ said, "Creators need to make efforts to protect data that should not be used in AI training." This was echoed by $R_{96}$: "If your product is a unique algorithm, then it is your responsibility to keep it closed" and by $R_{429}$: "If an individual or institution wishes to

have their data removed from a training dataset, they should have the right to request its removal."
(These respondents did not suggest how developers might accomplish such removal.)

In follow-up interviews, some respondents elaborated that they attempt to exercise such control
by limiting the types of projects for which they use GenAI tools, thereby restricting the tools'
access to their code. $RI_{118}$ specified that when choosing tools to use, "I'm not going to use a tool at
work that's going to expose my work to risk [...], whereas if it's just a personal project, I'm a lot
less concerned. I want to make sure I'm not going to violate anything major or there's nothing
truly weird like they're going to own my code."

> **Finding 15:** Some developers believed that the use of code as training data should be considered fair use by
> default and that the onus should be on developers to protect code that they do not wish to be included in
> training data.

## 5.3 AI-generated output that resembles existing code

Developers also had a range of views on whether it is wrongful for generated code to closely resemble
existing code. Some respondents compared what GenAI models do when producing code to what
humans have always done (manually) in creating derivative work. $R_{103}$ stated: "Being an open-source
software developer, I'm accustomed to seeing my codes copied or 'cloned' by other developers
without any notice, and I don't feel anything special from it." $R_{177}$ asserted, "All developers copy
code." $R_{265}$ noted, "What AI is doing now was being done by developers [throughout] the years
through forums, [S]tack [O]verflow, [G]ithub, etc." $R_{130}$ remarked, "Most developers had been using
Stack Overflow and OS[S] projects for years for similar purposes before AI became widely available."

> **Finding 16:** A group of developers did not see how GenAI technology was doing anything significantly
> different from traditional means of creating code, which reuses existing code created by others. The only
> distinction is the scale and ease by which it can be done.

For other respondents, their view on the wrongfulness of generated code that resembled code
they had developed depended on the nature and amount of the code being copied, an intuition
that tracks U.S. legal doctrine [8]. $R_{687}$ said that, in their view, wrongfulness depends on "if [the]
result is a simple hello world vs. word for word of an entire codebase [. . . ] In my current workflow,
expected outputs are only snippets, so I am indifferent about legality and attribution." This was an
attitude shared by other respondents. $R_{183}$ stated that "small pieces of code are rarely valuable in
themselves. For any given function there's only so many ways to write it. It's the entire thing that's
valuable." $R_{246}$ provided a concrete example, saying, "There are only so many ways you can create
a fetch request in Java[S]cript," and $R_{90}$ put it this way: "[T]o say setting up code to talk to a server
or render UI is unique is an exaggeration. Even your design, is what[,] a three column layout?
That is not someone's to own. It's all been done a million times already[,] and we are all reusing
founders technology going back to punch cards." $R_{573}$ summarized the concept: "Development is
typically a process of piecing together various snippets that are customized to a particular use
case; there is no expectation of 'ownership,' at least in the case of partial code blocks. While the
reproduction of an entire application/code base for use in generating revenue is problematic, the
sharing of common code is how the development community has progressed [until] now." And
$R_{445}$ offered a useful analogy: "[A]dding a license to a particular piece of code does not make sense.
You can add licenses to libraries and exact forms of algorithms. Adding a license to a piece of code
is like adding a copyright to individual strokes and colors in a painting."

> **Finding 17:** Some developers' views on the copying of code track U.S. copyright law: basic functions and
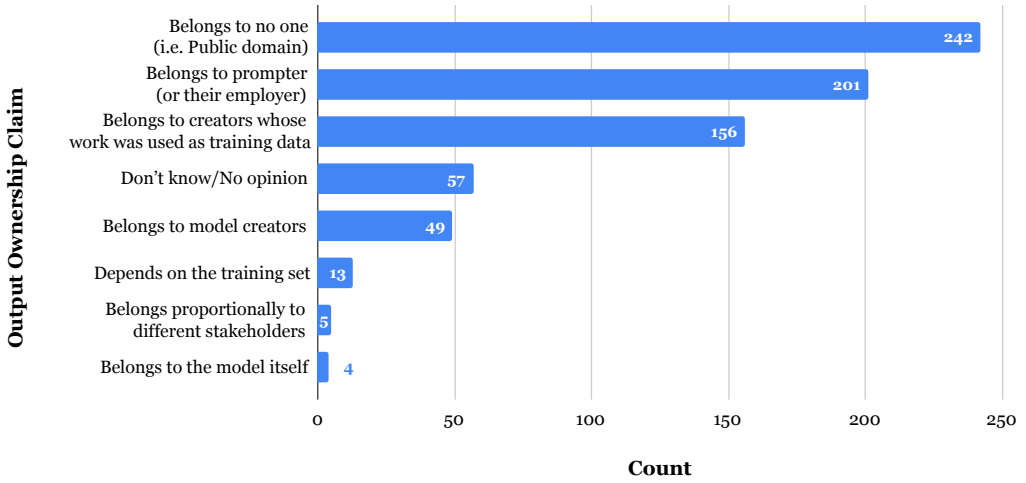> building blocks are not protectable and should be free for all to use.

Fig. 12. Developer perceptions on output ownership

## 5.4 On the ownership of AI-generated code

Our discussion to this point has focused on developers' views on potential liability for copying another's code, either in training data or in the generated output. But another important legal question is who would own any copyright to code created by GenAI models. In this subsection, we do not attempt to provide a legal analysis of this question; rather, our goal is to understand what developers think the answer to the question is.

Participants were presented with a list of potential options as to who owns the copyright to code generated through the use of AI and were invited to select all that they believed applied, with the opportunity to provide other options and to further explain their answers. Even though courts have yet to resolve this question, only 57 of 554 respondents familiar with GenAI tools (10.3%) indicated that they were not sure who owned the output or had no opinion on the matter. The rest of the respondents expressed an opinion and reported high confidence in their answers: 416 (75.1%) stated they were confident or even very confident.

*5.4.1 No one—works are in the public domain.* The most popular choice, selected by 242 participants (43.7%), was that AI-generated code belongs to no one. Some, including $R_7$, reasoned this way because they viewed the AI itself as the true author of the code: "The code belongs to the author—the AI model. It has no rights, so its work must be in the public domain." $R_{199}$ stated, similarly, "The generated code can't have any copyright as [it] wasn't made by a person." $R_{516}$ offered a summary: "Code generated through the use of AI is not copyrightable, because it fails the creative step required by law of being a product of the intellect. Instead it is produced by a non-sentient machine, lacks a creative step, an expression of the soul. It is therefore in the public domain."

Others, like $R_{140}$, selected this option to represent that they viewed the question as minimally important: "It's [going to] be ridiculously difficult to prove anything, and humanity's efforts should be spent on advancing technology and society, not fighting silly legal battles about who gets to call dibs on random binary stuff."

Several respondents invoked the inherently collaborative process of generating code through AI as a reason not to assign copyright ownership to only one individual or entity. $R_{74}$ gave an analogy:

"It'd be like saying who gets the copyright of what a child creates. [T]heir teachers? [T]heir parents? [T]heir friends? [T]hey all influence it in an unpredictable way. At least to my knowledge[,] we can't measure it." Putting it in more concrete terms, $R_{105}$ stated that the output of GenAI models is "a product of complex interactions between AI algorithms, training data, and user inputs, making traditional copyright assignment impractical. This perspective aligns with the open-source ethos, favoring communal access and innovation over individual ownership in the rapidly evolving field of AI and software development." $R_{518}$ said that the process of training and using LLMs "mirrors how all software development builds upon the foundation of previous programs and code [...], reflecting the collaborative nature of programming itself."

This sentiment regarding the nature of collaborative work also was echoed in responses, particularly from open-source developers, that seemed to suggest that even if the developer could claim copyright in AI output, they should refrain from asserting their rights in favor of the greater good. $R_{233}$ asserted that the focus should be on maximizing utility: "I've been contributing to open-source code for over a decade, consistently opting for the least restrictive licenses. I firmly believe in the value of sharing knowledge rather than imposing limitations. With the emergence of AI as a lasting tool, my stance is to focus on maximizing its utility for people rather than getting bogged down in minor details." $R_{546}$ said, similarly, "[W]hen I publish open-source code I desire that this code is re-used and helps as many people as possible. If AI models can scale this further and democratize the use of open-source code even more, I see this as highly desirable." $R_{695}$ concluded, "[I]f we're going to destroy a bunch of people's livelihoods, to produce something I find rather uncompelling, having it contribute to society as a whole via public domain would be great."

> **Finding 18:** The largest group of developers (242, 43.68%) thought that the output of GenAI belongs in the public domain. Some believe this because it is consistent with the collaborative nature of open-source software development, while others believe that, because AI models are not human, the output of those models cannot be subject to copyright law.

*5.4.2 The prompting developer or their employer.* The next most popular suggestion was that model output belonged to the developer that prompted the model or their employer, as selected by 201 respondents (36.3%).

Many participants who reached this conclusion believed that the use of code-generation tools is similar to existing processes that can result in copyrightable code. Respondents with this perspective, like $R_{173}$, characterized GenAI as "a tool like any other, used to augment the skills and direction of its user," where the assumption would be that the user owned the rights to whatever was created using the tool. $R_{22}$ provided an analogy: "[I]f a carpenter chooses to use a powerful automation tool like [a] CNC machine to make something, the final product is his (or his employer's), not the maker of the tool or [whoever] inspired the carpenter." $R_{236}$ summarized the sentiment as follows: "When one uses a tool to create something, you should own the result of your work." Ultimately, $R_{118}$ indicated, "If I can copyright code I learned by essentially copying what I learned from Stack Overflow, I should be able to copyright what I get from [...] Copilot. [...] If I can't, then the analysis should be similar. The focus on AI as the issue seems misguided. The only difference is scale."

For many developers, the prompt was an important starting point for the ownership question. $R_{129}$ stated plainly, "Without the prompt, there is no output." $R_{187}$ explained, "[C]ode generation requires prompt work, and the copyright on the code that is generated should be the result of that prompt work." For other developers, the claim to ownership relied on the assumption that code resulting from GenAI is not ready to use immediately upon generation but must be further refined by the developer. As explained by $R_{107}$, "AI code generation typically doesn't result in a perfectly working code. It's a sketch, or draft variant that will be modified and fine[-]tuned [...]. AI doesn't do development on its

own[;] it helps developers do their job." $R_{63}$ concluded that if the models were capable of "creating a large amount of code without [requiring] editing," copyright issues might be more complex, but $R_{63}$ was "not convinced that AI code generation is good enough for that." $R_{39}$ stated similarly, "The development process of using AI development tools usually involves AI generating an incorrect outline of one snippet of code in a larger system, that I then amend, correct and test. The additional work I put into it seems like it should constitute a new creative work that is copyrightable."

> **Finding 19:** The second largest group of developers (201, 36.3%) believed that the user of a GenAI tool should own the copyright to code generated with that tool because it is the user's contributions—creating the prompt and, typically, editing and further refining the result—that create the code. In other words, GenAI was seen as simply a new tool for creating content.

One challenge raised by some respondents, however, is that even if those supplying the prompts are deemed to be the owners of copyright in the resulting code, the nature of code creation with GenAI is such that similar code might result from two independently created prompts. Such activities could result, under U.S. law, in two separate copyrights, one for each generated work. $R_{98}$ found this legal conclusion to be counterintuitive: "[T]he developers who prompted can't also be the copyright owners – it's fairly simple to get very similar or identical outputs from [a] LLM by two different developers." $R_{276}$ expressed a similar concern, stating, "I believe that the prompt was the thing that caused the generation of the given code. And it could be generated again with a similar prompt," as did $R_{507}$, who wrote, "If I prompt the model to write code that I later sell, can I copyright it? Yes of course. [...] Can someone else prompt it the same way and maybe get the exact same code? Yes. Does that make sense? No. Which makes me think that AI copyright makes no sense."

Ultimately, some respondents noted, control over the resulting work incentivizes creative activity. $R_{300}$ asserted, "AI will not be a useful product to sell if the authors of the original training data or the model creators retain copyright on generated stuff, or if the output is uncopyrightable." $R_{114}$ mentioned that they "specifically use [c]ode [g]eneration tools that include in their ToS specific statements that copyright belongs to us (the customer)," and $R_{288}$ stated, "If the license does not end up with me at the end of the circus, I will use a different tool made by a different company trained off different data."

*5.4.3 The creators whose works make up the training data.* The next largest group, 156 respondents (28.2%), believed that rights in the generated code should belong to the individuals who created the code on which the GenAI model was trained. A common theme among such responses was a sense of fairness or justice, given GenAI's reliance on training data. As $R_{77}$ put it, "[I]f the code is extremely similar or identical to the training data[,] then the copyright should belong to the respective creator of this training data." $R_{13}$ wrote, "[I]f these humans had not written this code, this model would not exist," while $R_{25}$ echoed, "Without training data, AI can not do anything." Some responses expressed this view in stronger language, with $R_{31}$, for example, calling large language models "plagiarism engines."

Some respondents, by contrast, invoked the collaborative nature of development as a reason to deny copyright ownership claims by developers of the code used in the training data. $R_{157}$ noted, "[T]he neural networks require so much input for them to work that none of the input is really deterministic to the end work. [...] So the copyright can not belong to the people that just provided the training data. [If that were true], the copyright of [a] book could belong also to the authors of all the books that the author had read before." $R_{157}$ continued by providing another analogy from the creative arts: "[We don't give] copyright to the people that cleaned the theater and people that create the music[al] instruments and so on." Respondents also identified practical problems with

this ownership claim. $R_{187}$ said, "AI models are trained using [a] large [volume] of source code (some of which may be extremely similar) from different sources, so it is difficult to say which parts of the training material the generated code belongs to." In other words, an exact mapping from the training examples to the generated output might be difficult, if not impossible, to accomplish, particularly in situations with many similar or identical code instances (such as GitHub forks). As $R_{502}$ asked, "[H]ow can we be so sure whether that code was originally written by that user?"

> **Finding 20:** Some developers believed that fairness requires giving ownership of generated code to those developers who created the training data. However, others invoked the collaborative nature of development as a reason to deny copyright claims by those whose code was used in the training data and highlighted the practical problems with mapping training data to generated code in order to determine ownership.

*5.4.4 The creators of the model.* The least popular of the provided options, ranked even below the "don't know/no opinion" option, was that the creators of GenAI models should own the output, with only 49 participants (8.8%) indicating that they believed this to be the case. To justify this position, $R_1$ cited the "resources, time, money, effort and research that goes into creating novel systems." Similarly, $R_{441}$ believed that "the creators of the model should receive monetary compensation for training expenses."

By contrast, other respondents, such as $R_{306}$, indicated that while the ownership question is complicated, model creators were the one group with no real claim: "Again, it is difficult to be sure. What I do believe, however, is that it should never belong to the creators of the model, as they have no say on the output of the model, only its inner workings." $R_{98}$ expressed a similar sentiment: "The creators of the model are definitely not the copyright owners—LLMs are just text generators and as such can technically generate any text. You can force [a] LLM to 'write' literally any sequence of words you want—that [output] can't be copyrighted [by the model]." $R_{133}$ said "I can say very confidently that model creators should NOT get attribution for anything their model creates since they only aggregate the data into the model."

$R_{159}$ offered a pragmatic view: "If the code always belongs to the company that created and trained the model, there's no point in any company using it, especially freelance programmers and individuals." And $R_{157}$ suggested that if AI model creators held the copyright to generated code, that would mean for "any picture created by a camera, the copyright would also be held by the camera manufactur[er]. And we do not do this."

> **Finding 21:** The model vendors who develop GenAI models were the one group that the fewest participants (49) believed had an ownership claim to generated outputs.

*5.4.5 Other ownership claims.* Participants who selected multiple options conveyed that the answer to this question is often case-specific. Previous work that explores how changes to AI supply chains can influence the answers to copyright questions [75] supports this view. $R_{182}$ provided scenarios in which the ownership of generated code might vary: "Code generated by a model that I've prompted with a simple prompt (*e.g.*, 'Please show me a solution to the Fizz-Buzz challenge') should probably be public domain. However, code generated based on a prompt that includes significant copyright content (*e.g.*, 'Consider the style and techniques in the 50kB source file below, and then using those styles and relevant techniques, implement a new routine to [do something]. Before you start ask me any clarifying questions you think are important. [...]') could reasonably be protected as if it had been written by a human." $R_{152}$ wrote, "[It] depends on the jurisdiction, and on how much of the code was generated by the AI, and on how much the output was transformed by the prompter."

A few respondents chose multiple options because they believed that copyright ownership should be distributed among multiple groups. For example, $R_{191}$ called generated code "a co[-]creation

where we would have a percentage of ownership to apply. [...] The prompt, however, [would] belong entirely to [the] user." Similarly, $R_{264}$ pointed out that "Without original code, we couldn't train the AI tool; without the creator of the model, we will never have such a tool; without the [prompter], we couldn't get an ideal result. Everyone contributes to the AI generated results, so I think everyone should have some copyrights." $R_{527}$ suggested, "I would opt [for] [...] a shared benefits model, like [with] royalties in the entertainment industry, where the actual code owner gets a fair share, the model owner, and the end user who was capable of adding the right prompt into the model [...]." Such views were uncommon, and no respondent proposed a system for determining an appropriate allocation of ownership rights or the implications thereof.

> **Finding 22:** Some respondents believed that ownership depended on various factors, including the relevant jurisdiction, the amount of code generated, and how much the prompting developer contributed to the result. Others believed that multiple parties share claims to ownership and that a royalty-style system should be put into place.

Interestingly, a handful of respondents ($R_{86}$, $R_{489}$, $R_{593}$, $R_{715}$) advocated that the model itself should maintain rights to the content it generates, at least under certain conditions.

In the end, as $R_{54}$ noted, as GenAI becomes an increasingly common feature of all software, the question of who owns GenAI output may cease to matter: "[T]he issue [of ownership] becomes more clouded, though, when you consider the current situation, where most generative AI is done through a service[, n]ot directly with the model. Therefore the service that is providing the generated code or response can define in their own terms [of service] whether they own [the output] or not."

## 5.5 On the ownership and copying of prompts

Much of the current GenAI litigation, as of this writing, relates to the unauthorized use of preexisting work in training data or the copyright status of generated content. We were, however, also interested in understanding developers' views on the ownership of prompts, which, depending on their complexity, could be protected by copyright law. (As noted above, however, platforms' Terms of Service might govern users' rights in this regard.) About half of the 554 respondents familiar with GenAI tools (278, 50.2%) were indifferent as to whether a prompt they developed could be used by someone else without their permission, with 88 respondents (15.9%) indicating that they would be pleased if such use occurred and 63 (11.4%) extremely pleased. An additional 89 respondents (16.1%) would be displeased by the use of their prompts, and 36 (6.5%) extremely displeased. Overall, as shown in Figure 11, most developers in our survey stated that they would not care (or would even be pleased) if their prompts were reused by others, even without permission.

> **Finding 23:** Most respondents (77.4%) indicated that they would be indifferent or even pleased if their prompts were copied and used by others, even without their permission.

## 5.6 Perceived risk of potential copyright infringement

Ultimately, many developers perceived that the likelihood of legal enforcement resulting from improper usage of GenAI was low. $RI_{149}$ stated their belief: "Nobody's really looking, and if [copyrighted code] managed to get into a training data set somewhere, then [the copyright owner is] to blame." $RI_{39}$ stated in an interview: "If I technically am violating somebody's copyright for code that is never going to be pushed anywhere public, it's very implausible that I would be sued for that." $RI_{31}$ noted that they belonged to "a fairly small organization. We're not high on the radar. I don't think that we're at any legal risk of anything. If anybody's going to get sued, it's not going to

be us. It's going to be other people." $RI_{149}$ expressed a similar view: "[T]he likelihood of detection is so small that it's no longer a key concern."

> **Finding 24:** Individuals and smaller organizations may perceive less practical risk from using code from GenAI due to the low likelihood of any code that potentially infringes on another's copyright or violates another's license being discovered.

## 6 RQ₃: OTHER LEGAL CONCERNS THAT DEVELOPERS ANTICIPATE AS THE USE OF GENERATIVE AI INCREASES

**RQ₃** aimed to discover additional legal concerns (other than those relating to copyright) that developers anticipate as the use of GenAI increases. The possibility of liability may already be motivating organizations; $R_{34}$ indicated, "[My organization is] accounting for the possibility of legal risks and [has] an open source alternative if [O]penAI or [M]icrosoft are sued." Here we report these additional concerns that were on the minds of the surveyed developers. These include the aforementioned question of liability, the generation of non-code content, the creation of malicious content, risks involving data privacy, and the use of GenAI for software testing activities. A summary of key findings can be found in Table 7.

### 6.1 Liability concerns

Respondents mentioned concerns about possible tort liability related to the use of GenAI. Product liability was mentioned by 14 respondents (2.5%), and liability for code generated with bugs or defects was mentioned by 13 respondents (2.4%). Developers were concerned about who would be held accountable for security issues or other bugs introduced by AI-generated code. $R_0$ expressed that they "could see in the near future running [into] issues with security breaches caused by an AI code generator making sub-optimal choices. A contrived example[,] but imagine someone uses the AI to build a log in. The AI suggests that you store passwords in plain text in your database which is later leaked. The developer could argue the model creator (e.g. [O]penAI) was responsible for the breach."

Other respondents identified potential issues with AI agents, such as consumer-facing chatbots that purport to engage in contractual agreements with users [19, 87] or AI agents that appear to operate autonomously. $R_{149}$ said that "[f]unction calling"—using LLMs "to write and call functions"—was a potential concern (see examples in [14, 50]) and asked, "[W]hat happens if those called (which may have real world effects) cause harm—who is responsible?" $R_{204}$ raised the similar concern of rogue agents, which could be given a set of directives to autonomously carry out for malicious or nefarious purposes.

Many respondents explicitly mentioned that the use of generative technologies in important systems and domains posed a particularly high risk. $R_{109}$ asked "[i]f any medical, aviation, security or defen[s]e-based grade code is written by such systems, we better double check the work, and who takes the fall when such a system reaches production and fails?" $R_{245}$ commented more broadly: "Software liability is complicated, and although right now there seems to be a kind of stable equilibrium in terms of how liability is determined, AI may upend that equilibrium if we see individuals and companies with very large code bases that are mostly not authored by human beings. Some of these parties will inevitably end up trying to make the claim that liability should not attach in circumstances where errors such as 'negligence' were actually never carried [out] by a human being. This strikes me as a major concern for some types of software used in aerospace or medicine, etc."

Finally, 13 respondents (2.35%) were concerned that code might be generated without sufficient consideration of compliance issues, such as "ignoring GDPR [the EU's General Data Protection

Table 6. Main findings for RQ2: Developer perceptions of copyright issues related to GenAI

| | **5.1 Thoughts on Litigation** |
|---|---|
| F11 | While about half of participants anticipated that litigation would result in limitations on the use of GenAI, a vocal minority suggested that imposing such limitations would be difficult if not impossible, and others urged restraint, given the pace of technological developments. |
| | **5.2 Using one's own code in models' training data or as output** |
| | 5.2.1 Sentiments regarding the inclusion of one's own code in training data |
| F12 | Developers' views on having their code included in models' training data was context-dependent. If the code was open source, many developers were indifferent or even pleased to have their work included; this was not true for proprietary code. Developers also expressed concerns about the quality of code incorporated into the training data and whether increased use of GenAI tools would disincentivize developers from contribution to the OSS community. |
| | 5.2.2 Conditions around use of one's code in training data |
| F13 | Some respondents suggested that training data contributors should recieve monetary compensation or other benefits, such as free access to the trained model or free credits for the platform. However, forcing AI companies to provide monetary compensation to all contributors of the training data, while seemingly fair, would place an additional hurdle potentially blocking smaller AI companies from entering the market. |
| F14 | Respondents generally wanted licensing terms to be complied with by GenAI models, especially when producing code kidential to or very similar to work in the training set. |
| | 5.2.3 Fair use and replicated code |
| F15 | Some developers believed that the use of code as training data should be considered fair use by default and that the onus should be on developers to protect code that they do not wish to be included in training data. |
| | **5.3 AI-generated output that resembles existing code** |
| F16 | A group of developers did not see how GenAI technology was doing anything significantly different from traditional means of creating code, which reuses existing code created by others. The only distinction is the scale and ease at which it can be done. |
| F17 | Some developers' views on the copying of code track U.S. copyright law: basic functions and building blocks are not protectable and should be free for all to use. |
| | **5.4 On the ownership of AI-generated code** |
| | 5.4.1 No one-works are in the public domain |
| F18 | The largest group of developers (242, 43.68%) thought that the output of GenAI belongs in the public domain. Some believe this because it is consistent with the collaborative nature of open-source software development, while others believe that, because AI models are not human, the output of those models cannot be subject to copyright law. |
| | 5.4.2 The prompting developer or their employer |
| F19 | The second largerst group of developers (201, 36.3%) believed that a user of a GenAI tool should own the copyright to code generated with that tool because it is the user's contriubtions—creating the prompt and, typically, editing and further refining the result—that create the code. In other words, GenAI was seen as nothing more than a new tool for creating content. |
| | 5.4.3 The creators whose works make up the training data |
| F20 | Some developers believed that fairness requires giving ownership of generated code to those developers who created the training data. However, others invoked the collaborative nature of development as a reason to deny copyright claims by those whose code was used in the training data and highlighted the practical problems with mapping training data to generated code in order to determine ownership. |
| | 5.4.4 The creators of the model |
| F21 | The model vendors who developed GenAI models were the one group that the fewest participants (49) believed had an ownership claim to generated outputs. |
| | 5.4.5 Other ownership claims |
| F22 | Some respondents believed that ownership depended on various factors, including the relevant jurisdiction, the amount of code generated, and how much the prompting developer contributed to the result. Others believed that multiple parties share claims to ownership and that a royalty-style system should be put into place. |
| | **5.5 On the ownership and copying of prompts** |
| F23 | Most respondents (77.4%) indicated that they would be indifferent or even pleased if their prompts were copied and used by others, even without their permission. |
| | **5.6 Perceived risk of potential copyright infringement** |
| F24 | Individuals and smaller organizations may perceive less practical risk from using code from GenAI due to the low likelihood of any code that potentially infringes on another's copyright or violates another's license being discovered. |

Regulation]” ($R_{190}$). $R_{665}$ was concerned that “[i]f an AI generates specifications for a project, it may not consider existing laws and [so] generate software which does not comply with regulations.” Several respondents raised the related concern of using GenAI to draft legal documents such as software licenses, terms of service, and privacy policies, which could not only lead to an undesirable proliferation of such documents but also may result in language that does not have the desired legal effect. Similar problems of using AI to interpret or clarify legal standards were identified by $R_{527}$.

> **Finding 25:** Developers raised questions about who would be held accountable for bugs and other defects in GenAI outputs, especially for mission-critical systems such as medical, aerospace, and security. Developers were also concerned about who would be held responsible for actions that AI agents take autonomously.

## 6.2 Generating content other than code

Although our survey and interview questions focused primarily on copyright issues relating to code, 92 respondents (16.6%) also noted concerns about copyright issues relating to images and other types of content generated by AI, which might arise for developers in the context of “logo design [and] art generation for video games” ($R_{232}$) or fonts, icons, and images used in front-end web development, to take a few examples. (The legal status of outputs of image generation through AI tools such as MidJourney [91], StableDiffusion [115], and DALL-E [34] is currently a subject of debate [3, 10].)

In a similar vein, 23 respondents (4.2%) noted that tasks related to the development of user interfaces (UI) could lead to legal issues. $R_{54}$ asserted that “SVG images and logos are often created and defined directly with the code. Likewise design elements, CSS, and theme files are also involved in software design but are much easier targets for legal copyright as visual similarities are easier to point out.” $R_{30}$ added an additional concern: “I think relying on generative AI for UI/UX design can easily lead to essentially identical works.”

> **Finding 26:** The generation of other types of media (images, fonts, videos, audio, *etc.*) and theme/style-sheets can also pose potential legal challenges for development teams since such assets are often used in product creation.

## 6.3 Generation of malicious content

Seventeen survey respondents also highlighted the potential to use GenAI tools to produce a variety of malicous content, including malware or viruses, disinformation, illicit or unsavory material, and propaganda. While not a novel insight [41, 47, 103], developers are aware that the technical and skill barriers that have existed in creating malicious content have been lowered, enabling script-kiddies to carry out more sophisticated attacks with little to no domain knowledge. $R_{489}$ suggested that “AI [could be] used to maliciously attack someone (constant spam, emails, *etc.* too fast to fight back),” and $R_{411}$ mentioned that generative AI technologies could be used for carrying out search engine optimization (SEO) hacks. $R_{144}$ also raised the concern that models could generate malicious or toxic content that the end user did not solicit, which “could lead to legal repercussions in many EU countries.” While many mainstream LLMs have guardrails in place to prevent the generation of this type of content [101], unrestricted, open-weight models such as the Dolphin family of models [110] and alternatives such as WormGPT [95] can still be used by malicious actors, and jailbreaking techniques can be used to bypass safety rails even for more widely used models [81, 122].

> **Finding 27:** Respondents identified the purposeful generation of malicious content as a societal problem and legal concern going forward.

## 6.4 Data privacy, information leakage, and discovery of system vulnerabilities

Information leaks were mentioned by 33 (6%) respondents, and data privacy concerns by an additional 15 (2.7%). This concern has been explored in the existing literature [9, 22, 23], but it is of note that it is also a real world consideration for developers when choosing which GenAI tools to trust and use. Possible issues included proprietary code or information being produced by generative AI, leaked secrets such as API tokens and passwords, and the release of personally identifiable information. This risk results, as $R_{447}$ put it, from the fact that the "devs [...] didn't realize [sensitive information would] get hoovered up" in the training data. $R_{501}$ echoed this sentiment, noting that models were "[p]icking up private material that was committed accidentally."

> **Finding 28:** Developers and other stakeholders did not intend for or anticipate that their data would be collected *en masse*, so developers are concerned that private and sensitive information can easily make its way into training sets.

Similarly, there were concerns that generative models could be used by malicious actors for hacking (6), reverse engineering (6), and jailbreaking (1) systems. $R_{393}$ said that people could use "AI to hack/discover security vulnerabilities [...]." $R_{139}$ said that people might gain "access to secure systems through intelligent vulnerability discovery and exploitation." $R_{352}$ worried that AI tooling "could make it easier to jailbreak secure applications by learning salts and similar common practices in infrastructure to breach a site or app."

> **Finding 29:** Generative AI, particularly LLMs, can not only leak information that was contained in the original training data but can also be used by malicious actors to exploit systems.

Respondents were also cognizant of the dangers of including private or proprietary information in prompts, fearing that those same prompts would be used in training sets of future models. $RI_{31}$ explained, "Because I'm almost only doing work stuff, and in the enterprise environment, I can't run the risk of the training set pulling security keys and that kind of thing. Because it's been demonstrated that you can trick those things into outputting data from the training set." $RI_{516}$ said, similarly, that "I use [GenAI] only for things that are very generic and do not use code that would not already go in the open."

## 6.5 Testing and test generation

Tasks related to testing, such as generating test cases, were identified by 28 (5.1%) respondents as having potential legal ramifications. Some respondents identified software testing itself as a shortcoming of generative AI's capabilities; as $R_{478}$ put it, "[W]ho will test the generated test?" Beyond this, participants asserted that the application of generative AI for testing could lead to legal problems in several ways, including use for penetration testing ($R_{299}$) and security audits ($R_{364}$). The specter of liability also emerged ($R_{314}$): "AI models can automate testing processes, but if automated testing leads to false negatives or positives that result in financial or reputational harm, it may result in legal disputes." $R_{434}$ also cautioned that "[g]enerating random datasets (csv, image, audio, video, etc.) for software testing may violate someone's personal rights."

> **Finding 30:** Developers had concerns related to AI-generated tests, including their completeness and reliability.

Table 7. Main findings for RQ3: Other legal concerns anticipated by developers

| | **6.5 Liability concerns** |
|---|---|
| F25 | Developers raised questions about who would be held accountable for bugs and other defects in GenAI outputs, especially for mission-critical systems such as medical, aerospace, and security. Developers were also concerned about how would be held responsible for actions that AI agents take autonomously. |
| | **6.1 Generating content other than code** |
| F26 | The generation of other types of media (images, fonts, videos, audio, *etc.* and theme/style-sheets can also pose potential legal challenges for development teams since such assets are often used in product creation. |
| | **6.2 Generation of malicious content** |
| F27 | Respondents identified the purposeful generation of malicious content as a societal problem and legal concern going forward. |
| | **6.3 Data privacy and information leakage** |
| F28 | Developers and other stakeholders did not intend for or anticipate that their data would be collected *en masse*, so developers are concerned that private and sensitive information can easily makes its way into training sets. |
| F29 | Generative AI, particularly LLMs, can not only leak information that was contained in the original training data but can also be used by malicious actors to exploit systems. |
| | **6.4 Testing and test generation** |
| F30 | Developers had concerns related to AI-generated tests, including their completeness and reliability. |

## 7 THREATS TO VALIDITY

### 7.1 Construct validity

Given the design of our study, the results represent developers' communicated perceptions, which may or may not reflect their or their organizations' actual practices. We have avoided making any claim connecting such perceptions to actual practices. While we sought to include a wide range of perspectives, the demographic information included in this study was self-reported by respondents and was not independently verified. We also acknowledge that neither the survey respondents nor the interviewees may constitute a representative sample of developers. However, through our coding process, we aimed to report those views and perceptions that were held by several developers.

### 7.2 Internal validity

We followed best practices when developing the questionnaire to ensure that the questions were clear and would not bias the participants toward a particular answer. Additionally, we also conducted a small pilot study with graduate students in our research lab to check readability, length, and content. Given that participation in the survey was done on an entirely voluntary basis from a single source of participants, we are aware of the problems introduced by self-selection bias. However, the number of participants gives us confidence that the results capture various developers' views. In the qualitative analysis of survey responses, we mitigated confirmation bias by having multiple annotators examine the data independently and then reconvene to discuss and resolve any disagreements, rigorously following the best practices of qualitative analysis. Our conclusions were derived from data analysis.

### 7.3 External validity

To investigate developers' perspectives, we rely on GitHub developers who showed interest in repositories of GenAI tools for coding. While our respondents may represent developers working in different domains and organizations, we are also aware that they may overrepresent views typical of open-source ecosystems and underrepresent those of closed-source developers. Although we do not claim generalizability, many respondents work for companies or in education, which may lead to perspectives that vary from those in the open-source community. Moreover, while we had participants from 73 countries spanning 6 continents, nearly a third of respondents came from the U.S., so our results may apply most directly to the U.S. context and legislation.

# 8  DISCUSSION, IMPLICATIONS, AND CONCLUSIONS

## 8.1  Developing new licenses for software used in AI training

Open-source software (OSS) is used to train GenAI models, and so licensing is an important consideration. Respondents expressed a range of views on the permissibility of using OSS as training data (Section 5.2.1), with some condoning such use (Section 5.2.3) and others believing that OSS developers should take proactive steps if they want their work to be excluded from such use, such as including restrictions in a license (Section 5.2.2). For many years, there has been a standard set of licenses in the OSS community, and drafting new licenses ("license proliferation") has been discouraged, as it makes compliance tasks more difficult [16, 124]. However, the text of the most popular OSS licenses was drafted decades ago and did not anticipate technological developments such as GenAI. This means that it may be unclear whether or how the terms of existing OSS licenses apply to scenarios in which OSS projects are used to train, develop, and distribute AI systems.

In the past, updated versions of popular licenses have been drafted to account for new technology [116, 120], and AI model-specific licenses (*e.g.*, OpenRAIL [109]) have emerged that impose ethical usage requirements. Thus, the new legal challenges introduced by large-scale training and GenAI may necessitate the drafting of new OSS licenses or the revision of existing licenses so as to provide guidance on these issues. Importantly, since most existing OSS licenses were drafted by developers, who typically do not have legal experience, their terms are sometimes difficult to align with legal doctrine. This has given rise to the less-than-ideal situation in which drafting organizations issue FAQs explaining the intent of their licenses [124]. If those with legal expertise, preferably with accompanying SE backgrounds, participate in drafting these new licenses, it is possible that ambiguous and unclear language can be avoided to the largest extent possible.

## 8.2  Disclosing training data and tracking data provenance

Participants expressed a desire that information about the content and provenance of training data be made available for transparency (Sections 5.2.2 and 5.2.3), specifically to fulfill attribution requirements, to allow community checks for security risks or harmful content, to provide a mechanism for creatives to determine if their works have been included, to give users information through which to evaluate tools on ethical and philosophical grounds, and to allow organizations to ascertain legal risk from using tools trained on datasets.

Providing complete information on the content in training material comes with both technical and legal hurdles. It may appear sensible at first to apply proposed DataBOMs (Data Bills of Materials) [17, 117] to the problem, but the issue becomes one of scale. Unlike more canonical machine learning techniques, LLMs are trained on billions, if not trillions, of data points collected from large swaths of the internet. Such datasets are hard to navigate and even harder to curate [76], and any attempt to provide a complete list of the training data must consider that a work might appear in several different sources. For example, even if a creator's original work is not directly included in the dataset, it might appear in a StackOverflow post that is included in the dataset. As such, a single DataBOM will likely be insufficient.

Data provenance is similarly complicated. While such functionality has been achieved by newer models, like ChatGPT-4 and Gemini, this relies on Retrieval Augmented Generation (RAG)-style systems [84]. Achieving correct citations for model generations that do not rely on external data sources accessible through RAG or provided in a prompt is still an open research area.

Since new LLM models can be created by fine-tuning existing base models, it also becomes necessary to appropriately track the provenance and lineage of models. Failure to do so could result in not all training data being known or disclosed. These relationships emerging from the

AI supply chain could be tracked using standardized model cards [92] or AIBOMs (AI Bills of Materials) [117, 127].

Additionally, entities might train or fine-tune a model using proprietary or otherwise sensitive data sources. This creates a tension between transparency and protecting intellectual property or proprietary secrets. Should entities be required to provide traceability and transparency information for the models they release or otherwise make available? Balancing provenance with these data privacy concerns will also be a challenge for future work.

Finally, some respondents called for a right to be removed from training data. Apart from the difficulty of determining whether a data point is included in a training dataset, the further question of what it means to be "removed" and how such removal can be accomplished effectively remain unanswered [75, 130].

### 8.3 New AI advancements will likely amplify legal challenges

Many of the study participants expressed the sentiment that "tools are tools" and, as such, developers should hold the copyright to GenAI output (Section 5.4.2). Much of this understanding is based on the current limitations of existing tools, which, in many cases, require developers to manually correct, edit, or adapt the solution provided by the AI before the code can be used. A few respondents, however, indicated that they anticipated more complicated legal challenges ahead once the technology advances and can produce entire software packages composed of multiple files that largely resemble existing licensed components, a process already under development [49, 125]. The question will then be whether developers' attitudes about generated code ownership will change based on the nature and scope of developers' own contributions. Likewise, while developers' intuitions about code usage may currently align, at least to some extent, with U.S. copyright law, there may well be a future misalignment as both the technology and the law evolve. Having a better understanding of developer expectations and perceptions could help guide regulators and decision-makers in this area.

### 8.4 Concluding remarks

Although there are still many open questions regarding the use of generative AI and copyright, the biggest priority for future research is to address problems related to data provenance. Developers expect tools that can provide them with citations and will not lead them unwittingly into legal trouble. Likewise, they need tools that allow them to document and track code generated by LLMs for transparency and control of legal risk. Given that this paper primarily examines these issues from a perspective grounded in U.S. law, we encourage further research to interpret our results against regulatory efforts elsewhere.

## REFERENCES

[1] 2023. Letter from Robert J. Kasunic, Associate Register of Copyrights and Director of Registration Policy and Practice, U.S. Copyright Office, to Van Lindberg, Taylor English Duma LLP (February 21, 2023) , https://www.copyright.gov/docs/zarya-of-the-dawn.pdf.

[2] 2023. Andy Warhol Foundation for the Visual Arts, Inc. v. Goldsmith, 143 S. Ct. 1258 (2023).

[3] 2023. Getty Images (US) Inc v Stability AI Inc. 1:23-cv-00135, (D. Del.).

[4] 2024. Online replication package. https://anonymous.4open.science/r/ai_legal_survey-37D0/.

[5] [n.d.]. Qualtrics. https://www.qualtrics.com/. Accessed: 2023-21-06.

[6] 17 U.S.C. § 102(a) [n. d.]. 17 U.S.C. § 102(a).

[7] 17 U.S.C. § 106 [n. d.]. 17 U.S.C. § 106.

[8] 17 U.S.C. §102(b) [n. d.]. 17 U.S.C. §102(b).

[9] John Abascal, Stanley Wu, Alina Oprea, and Jonathan Ullman. 2023. TMI! Finetuned Models Leak Private Information from their Pretraining Data. *Proc. Priv. Enhancing Technol.* 2024 (2023), 202–223. https://api.semanticscholar.org/CorpusID:259064156

[10] Andersen v. Stability AI Ltd. [n. d.]. Andersen v. Stability AI Ltd. Andersen v. Stability AI Ltd., 3:23-cv-00201, (N.D. Cal.).

[11] Anthropic Models [n. d.]. Models. https://docs.anthropic.com/en/docs/about-claude/models. Accessed: 2024-17-10.

[12] Authors Guild, Inc. v. Google, Inc. 2015. Authors Guild, Inc. v. Google, Inc., 804 F.3d 202 (2d Cir. 2015)..

[13] Authors Guild v. OpenAI Inc. 2023. Authors Guild v. OpenAI Inc. Authors Guild v. OpenAI Inc., 1:23-cv-08292, (S.D.N.Y.).

[14] AutoGPT [n. d.]. AutoGPT. https://github.com/Significant-Gravitas/AutoGPT.

[15] B. Rozière *et al.*. 2024. Code Llama: Open Foundation Models for Code. arXiv:2308.12950 [cs.CL]

[16] Mahak Bandi. 2019. All About Open Source Licenses. https://fossa.com/blog/what-do-open-source-licenses-even-mean/. Accessed: 2023-24-09.

[17] Iain Barclay, Alun Preece, Ian Taylor, and Dinesh Verma. 2019. Towards Traceability in Data Ecosystems using a Bill of Materials Model. *arXiv preprint arXiv:1904.04253* (2019).

[18] Shraddha Barke, Michael B James, and Nadia Polikarpova. 2023. Grounded Copilot: How Programmers Interact with Code-Generating Models. *Proceedings of the ACM on Programming Languages* 7, OOPSLA1 (2023), 85–111.

[19] Ashley Belanger. 2024. Air Canada Has to Honor a Refund Policy Its Chatbot Made Up. https://www.wired.com/story/air-canada-chatbot-refund-policy.

[20] Blake Brittain. 2023. Authors sue Meta, Microsoft, Bloomberg in latest AI copyright clash. https://www.reuters.com/legal/litigation/authors-sue-meta-microsoft-bloomberg-latest-ai-copyright-clash-2023-10-18/.

[21] california-regulation 2024. Governor Newsom announces new initiatives to advance safe and responsible AI, protect Californians. https://www.gov.ca.gov/2024/09/29/governor-newsom-announces-new-initiatives-to-advance-safe-and-responsible-ai-protect-californians/.

[22] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021. Extracting Training Data from Large Language Models. In *30th USENIX Security Symposium (USENIX Security 21)*. 2633–2650.

[23] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2020. Extracting Training Data from Large Language Models. In *USENIX Security Symposium*. https://api.semanticscholar.org/CorpusID:229156229

[24] codeium [n. d.]. codeium.vim. https://github.com/Exafunction/codeium.vim.

[25] cody [n. d.]. cody. https://github.com/sourcegraph/cody.

[26] colorado-bill 2024. Consumer Protections for Artificial Intelligence. https://leg.colorado.gov/bills/sb24-205.

[27] European Commission. 2024. The EU copyright legislation. https://digital-strategy.ec.europa.eu/en/policies/copyright-legislation.

[28] Kattiana Constantino, Mauricio Souza, Shurui Zhou, Eduardo Figueiredo, and Christian Kästner. 2023. Perceptions of open-source software developers on collaborations: An interview and survey study. *Journal of Software: Evolution and Process* 35, 5 (2023), e2393.

[29] CopyrightCatcher [n. d.]. Introducing CopyrightCatcher, the first Copyright Detection API for LLMs. Accessed: Mar. 22 2024. https://www.patronus.ai/blog/introducing-copyright-catcher.

[30] Council of the European Union. 2024. Article 53: Obligations for Providers of General-Purpose AI Models. https://artificialintelligenceact.eu/article/53/.

[31] Council of the European Union. 2024. Proposal for a Regulation of the European Parliament and of the Council laying down harmonised rules on artificial intelligence (Artificial Intelligence Act) and amending certain Union legislative acts. https://digital-strategy.ec.europa.eu/en/policies/regulatory-framework-ai.

[32] cptX [n. d.]. cptX. https://github.com/maxim-saplin/cptX.

[33] Carys J Craig. 2024. THE AI-COPYRIGHT TRAP. *Available at SSRN* (2024).

[34] DALL-E 3 [n. d.]. DALL-E 3. https://openai.com/dall-e-3.

[35] DOE 1 et al v. GitHub [n. d.]. DOE 1 et al v. GitHub. 4:22-cv-06823, (N.D. Cal.).

[36] Michael Duan, Anshuman Suri, Niloofar Mireshghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. 2024. Do Membership Inference Attacks Work on Large Language Models? *arXiv preprint arXiv:2402.07841* (2024).

[37] Christof Ebert and Panos Louridas. 2023. Generative AI for software practitioners. *IEEE Software* 40, 4 (2023), 30–38.

[38] Emilio Ferrara. 2023. Fairness and Bias in Artificial Intelligence: A Brief Survey of Sources, Impacts, and Mitigation Strategies . *Sci* 6, 1 (2023), 3.

[39] figma [n. d.]. Figma. https://www.figma.com/.

[40] Anna Filippova and Hichang Cho. 2016. The effects and antecedents of conflict in free and open source software development. In *Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing*. 705–716.

[41] Lothar Fritsch, Aws Jaber, and Anis Yazidi. 2022. An overview of artificial intelligence used in malware. In *Symposium of the Norwegian AI Society*. Springer, 41–51.

[42] Gemini Team et al. 2023. Gemini: A Family of Highly Capable Multimodal Models. arXiv:2312.11805 [cs.CL]

[43] Generative Artificial Intelligence and Copyright Law [n. d.]. Generative Artificial Intelligence and Copyright Law. https://crsreports.congress.gov/product/pdf/LSB/LSB10922.

[44] Daniel J Gervais, Noam Shemtov, HARALAMBOS MARMANIS, and CATHERINE ZALLER ROWLAND. 2024. The Heart of the Matter: Copyright, AI Training, and LLMs. (2024).

[45] GitHub. [n. d.]. GitHub Copilot. Accessed: Mar 22 2024. https://copilot.github.com.

[46] GitHub REST API documentation [n. d.]. GitHub REST API documentation. https://docs.github.com/en/rest.

[47] Josh A Goldstein, Jason Chao, Shelby Grossman, Alex Stamos, and Michael Tomz. 2024. How persuasive is AI-generated propaganda? *PNAS nexus* 3, 2 (2024), pgae034.

[48] Google LLC v. Oracle America, Inc. [n. d.]. Google LLC v. Oracle America, Inc., 141 S. Ct. 1183 (2021)..

[49] gpt-pilot [n. d.]. gpt-pilot. https://github.com/Pythagora-io/gpt-pilot.

[50] gptengineer [n. d.]. gpt-engineer. https://github.com/gpt-engineer-org/gpt-engineer.

[51] grammarly [n. d.]. Transforming How the World Communicates Through AI. https://www.grammarly.com/ai.

[52] Robert M. Groves, Floyd J. Jr. Fowler, Mick P. Couyper, James M. Lepkowski, Eleanor Singer, and Roger Tourangeau. 2009. *Survey Methodology, 2nd edition*. Wiley.

[53] Andres Guadamuz. 2024. A Scanner Darkly: Copyright Liability and Exceptions in Artificial Intelligence Inputs and Outputs. *GRUR International* 73, 2 (2024), 111–127.

[54] Maanak Gupta, CharanKumar Akiri, Kshitiz Aryal, Eli Parker, and Lopamudra Praharaj. 2023. From ChatGPT to ThreatGPT: Impact of Generative AI in Cybersecurity and Privacy. *IEEE Access* (2023).

[55] Runzhi He, Hao He, Yuxia Zhang, and Minghui Zhou. 2023. Automating dependency updates in practice: An exploratory study on github dependabot. *IEEE Transactions on Software Engineering* 49, 8 (2023), 4004–4022.

[56] Peter Henderson, Jieru Hu, Mona Diab, and Joelle Pineau. 2024. Rethinking Machine Learning Benchmarks in the Context of Professional Codes of Conduct. In *Proceedings of the Symposium on Computer Science and Law*. 109–120.

[57] Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A Lemley, and Percy Liang. 2023. Foundation Models and Fair Use. *arXiv preprint arXiv:2303.15715* (2023).

[58] Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352* (2023).

[59] Stephen Hood. 2023. llamafile: bringing LLMs to the people, and to your own computer. https://future.mozilla.org/builders/news_insights/introducing-llamafile/. Accessed: 2024-09-11.

[60] Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. 2022. Membership Inference Attacks on Machine Learning: A Survey. *ACM Computing Surveys (CSUR)* 54, 11s (2022), 1–37.

[61] Yu Huang, Denae Ford, and Thomas Zimmermann. 2021. Leaving my fingerprints: Motivations and challenges of contributing to OSS for social good. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 1020–1032.

[62] HuggingChat [n. d.]. HuggingChat. https://huggingface.co/chat/. Accessed: 2024-08-11.

[63] Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A Choquette-Choo, and Nicholas Carlini. 2023. Preventing Generation of Verbatim Memorization in Language Models Gives a False Sense of Privacy. In *Proceedings of the 16th International Natural Language Generation Conference*. Association for Computational Linguistics, 28–53.

[64] Jing Jiang, David Lo, Xinyu Ma, Fuli Feng, and Li Zhang. 2017. Understanding inactive yet available assignees in GitHub. *Information and Software Technology* 91 (2017), 44–55.

[65] Mitchell Joblin, Sven Apel, Claus Hunsen, and Wolfgang Mauerer. 2017. Classifying developers into core and peripheral: An empirical study on count and network metrics. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*. IEEE, 164–174.

[66] Theodoros Karathanasis. 2023. EU Copyright Directive: A 'Nightmare' for Generative AI Researchers and Developers? https://ai-regulation.com/eu-copyright-directive-a-nightmare-for-gai/.

[67] Barbara A. Kitchenham and Shari Lawrence Pfleeger. 2002. Principles of Survey Research Part 2: Designing a Survey. *ACM SIGSOFT Software Engineering Notes* 27, 1 (2002), 18–20.

[68] Barbara A. Kitchenham and Shari Lawrence Pfleeger. 2002. Principles of Survey Research: Part 3: Constructing a Survey Instrument. *ACM SIGSOFT Software Engineering Notes* 27, 2 (2002), 20–24.

[69] Barbara A. Kitchenham and Shari Lawrence Pfleeger. 2002. Principles of Survey Research Part 4: Questionnaire Evaluation. *ACM SIGSOFT Software Engineering Notes* 27, 3 (2002), 20–23.

[70] Barbara A. Kitchenham and Shari Lawrence Pfleeger. 2002. Principles of Survey Research: Part 5: Populations and Samples. *ACM SIGSOFT Software Engineering Notes* 27, 5 (2002), 17–20.

Trevor Stalnaker, Nathan Wintersgill, Oscar Chaparro, Laura A. Heymann, Massimiliano Di Penta, Daniel M German, and Denys Poshyvanyk

40

[71] Barbara A. Kitchenham and Shari Lawrence Pfleeger. 2003. Principles of Survey Research Part 6: Data Analysis. *ACM SIGSOFT Software Engineering Notes* 28, 2 (2003), 24–27.

[72] Paul Krill. [n. d.]. GitHub survey finds nearly all developers using AI coding tools. https://www.infoworld.com/article/3489925/github-survey-finds-nearly-all-developers-using-ai-coding-tools.html.

[73] Logan Kugler. 2024. Who Owns AI's Output? *Commun. ACM* (2024). https://api.semanticscholar.org/CorpusID:273132452

[74] Jose Antonio Lanz. [n. d.]. AI Art Wars: Japan Says AI Model Training Doesn't Violate Copyright. https://finance.yahoo.com/news/ai-art-wars-japan-says-185350499.html.

[75] Katherine Lee, A Feder Cooper, and James Grimmelmann. 2023. Talkin' 'Bout AI Generation: Copyright and the Generative-AI Supply Chain. *arXiv preprint arXiv:2309.08133* (2023).

[76] Katherine Lee, A. Feder Cooper, James Grimmelmann Grimmelmann, and Daphne Ippolito. 2023. The Devil is in the Training Data. https://genlaw.org/explainers/training-data.html.

[77] Mark A Lemley and Bryan Casey. 2020. Fair learning. *Tex. L. Rev.* 99 (2020), 743–785.

[78] Jiawei Li and Iftekhar Ahmed. 2023. Commit message matters: Investigating impact and evolution of commit message quality. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 806–817.

[79] Jenny T Liang, Chenyang Yang, and Brad A Myers. 2024. A large-scale survey on the usability of AI programming assistants: Successes and challenges. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*. 1–13.

[80] Jenny T Liang, Thomas Zimmermann, and Denae Ford. 2022. Understanding skills for OSS communities on GitHub. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 170–182.

[81] Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, and Kailong Wang. 2024. A Hitchhiker's Guide to Jailbreaking ChatGPT via Prompt Engineering. In *Proceedings of the 4th International Workshop on Software Engineering and AI for Data Quality in Cyber-Physical Systems/Internet of Things*. 12–21.

[82] Llama.cpp [n. d.]. Llama.cpp. https://github.com/ggerganov/llama.cpp. Accessed: 2024-08-11.

[83] Nicola Lucchi. 2023. ChatGPT: A Case Study on Copyright Challenges for Generative Artificial Intelligence Systems. *European Journal of Risk Regulation* (2023), 1–23.

[84] Lijia Ma, Xingchen Xu, and Yong Tan. 2024. Crafting Knowledge: Exploring the Creative Mechanisms of Chat-Based Search Engines. *arXiv preprint arXiv:2402.19421* (2024).

[85] Vahid Majdinasab, Michael Joshua Bishop, Shawn Rasheed, Arghavan Moradidakhel, Amjed Tahir, and Foutse Khomh. 2024. Assessing the Security of GitHub Copilot's Generated Code-A Targeted Replication Study. In *2024 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 435–444.

[86] Vahid Majdinasab, Amin Nikanjam, and Foutse Khomh. 2024. Trained Without My Consent: Detecting Code Inclusion In Language Models Trained on Code. *arXiv preprint arXiv:2402.09299* (2024).

[87] Tom Malley. 2023. AI HAVE A DEAL Driver uses ChatGPT hack to get dealer to agree to sell new car for $1 in 'legally binding deal' in blow for AI rollout. https://www.the-sun.com/motors/9888857/driver-uses-ai-loophole-buy-new-car-1.

[88] Eira May. [n. d.]. Where developers feel AI coding tools are working—and where they're missing the mark. https://stackoverflow.blog/2024/09/23/where-developers-feel-ai-coding-tools-are-working-and-where-they-re-missing-the-mark/.

[89] Shiona McCallum. 2023. ChatGPT banned in Italy over privacy concerns. https://www.bbc.com/news/technology-65139406.

[90] Meta. 2024. Introducing Meta Llama 3: The most capable openly available LLM to date. https://ai.meta.com/blog/meta-llama-3/.

[91] Midjourney [n. d.]. Midjourney. https://www.midjourney.com/home.

[92] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model Cards for Model Reporting. In *Proceedings of the conference on fairness, accountability, and transparency*. 220–229.

[93] Joao Pedro Moraes, Ivanilton Polato, Igor Wiese, Filipe Saraiva, and Gustavo Pinto. 2021. From one to hundreds: multi-licensing in the JavaScript ecosystem. *Empirical Software Engineering* 26 (2021), 1–29.

[94] Seth Neel and Peter Chang. 2023. Privacy Issues in Large Language Models: A Survey. *arXiv preprint arXiv:2312.06717* (2023).

[95] Newsroom. 2023. WormGPT: New AI Tool Allows Cybercriminals to Launch Sophisticated Cyber Attacks. https://thehackernews.com/2023/07/wormgpt-new-ai-tool-allows.html.

[96] United States Copyright Office. 2023. Copyright Registration Guidance: Works Containing Material Generated by Artificial Intelligence. 16190 Federal Register, Vol. 88, No. 51.

[97] United States Copyright Office. 2024. Artificial Intelligence Study. https://www.copyright.gov/policy/artificial-intelligence/.

[98] Ollama [n. d.]. Ollama. https://ollama.com/. Accessed: 2024-08-11.

[99] Oobabooga [n. d.]. Oobabooga Text Generation WebUI. https://github.com/oobabooga/text-generation-webui. Accessed: 2024-08-11.

[100] OpenAI. [n. d.]. ChatGPT https://openai.com/blog/chatgpt. Last accessed: March 2024.

[101] OpenAI. 2024. OpenAI safety update. https://openai.com/index/openai-safety-update/.

[102] Abraham Naftali Oppenheim. 2000. *Questionnaire design, interviewing and attitude measurement.* Bloomsbury Publishing.

[103] Yin Minn Pa Pa, Shunsuke Tanizaki, Tetsui Kou, Michel Van Eeten, Katsunari Yoshioka, and Tsutomu Matsumoto. 2023. An attacker's dream? exploring the capabilities of chatgpt for developing malware. In *Proceedings of the 16th Cyber Security Experimentation and Test Workshop*. 10–18.

[104] Perplexity [n. d.]. What is Perplexity? https://www.perplexity.ai/hub/faq/what-is-perplexity. Accessed: 2024-08-11.

[105] Shari Lawrence Pfleeger and Barbara A. Kitchenham. 2001. Principles of Survey Research: Part 1: Turning Lemons into Lemonade. *ACM SIGSOFT Software Engineering Notes* 26, 6 (2001), 16–18.

[106] R. Li *et al.*. 2023. StarCoder: may the source be with you. *arXiv preprint arXiv:2305.06161* (2023).

[107] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*. PMLR, 28492–28518.

[108] Asha Rajbhoj, Akanksha Somase, Piyush Kulkarni, and Vinay Kulkarni. 2024. Accelerating Software Development Using Generative AI: ChatGPT Case Study. In *Proceedings of the 17th Innovations in Software Engineering Conference*. 1–11.

[109] Responsible AI. 2022. Big Science Open Rail-M License https://www.licenses.ai/blog/2022/8/26/bigscience-open-rail-m-license.

[110] Aresh Sarkari. 2024. Exploring Uncensored LLM Model – Dolphin 2.9 on Llama-3-8b. https://askaresh.com/2024/05/02/exploring-uncensored-llm-model-dolphin-2-9-on-llama-3-8b/.

[111] Sk Golam Saroar and Maleknaz Nayebi. 2023. Developers' perception of GitHub Actions: A survey analysis. In *Proceedings of the 27th International Conference on Evaluation and Assessment in Software Engineering*. 121–130.

[112] Agnia Sergeyuk, Yaroslav Golubev, Timofey Bryksin, and Iftekhar Ahmed. 2024. Using AI-Based Coding Assistants in Practice: State of Affairs, Perceptions, and Ways Forward. *CoRR* abs/2406.07765 (2024).

[113] The Yomieru Shinbun. June 10, 2023. Intellectual Property Plan Signals Reversal on AI Policy. Japan News.

[114] Donna Spencer. 2009. *Card sorting: Designing usable categories.* Rosenfeld Media.

[115] Stability AI [n. d.]. Stability AI. https://stability.ai/stable-image.

[116] Richard Stallman. [n. d.]. Why Upgrade to GPLv3. https://www.gnu.org/licenses/rms-why-gplv3.

[117] Trevor Stalnaker, Nathan Wintersgill, Oscar Chaparro, Massimiliano Di Penta, Daniel M German, and Denys Poshyvanyk. 2024. BOMs Away! Inside the Minds of Stakeholders: A Comprehensive Study of Bills of Materials for Software Systems. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*. 1–13.

[118] tabby [n. d.]. tabby. https://github.com/TabbyML/tabby.

[119] Tabnine [n. d.]. The AI code assistant you control. https://www.tabnine.com/. Accessed: 2024-25-10.

[120] FOSSA Editorial Team. 2021. Open Source Software Licenses 101: The AGPL License. https://fossa.com/blog/open-source-software-licenses-101-agpl-license/.

[121] The New York Times Company v. Microsoft Corporation [n. d.]. The New York Times Company v. Microsoft Corporation, No. 1:23-cv-11195 (S.D.N.Y., filed Dec. 27, 2023),https://nytco-assets.nytimes.com/2023/12/NYT_Complaint_Dec2023.pdf.

[122] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2024. Jailbroken: How does llm safety training fail? *Advances in Neural Information Processing Systems* 36 (2024).

[123] Whisper [n. d.]. Whisper. https://github.com/openai/whisper. Accessed: 2024-11-12.

[124] Nathan Wintersgill, Trevor Stalnaker, Laura A Heymann, Oscar Chaparro, and Denys Poshyvanyk. 2024. "The Law Doesn't Work Like a Computer": Exploring Software Licensing Issues Faced by Legal Practitioners. In *Proceedings of the ACM on Software Engineering*, Vol. 1. ACM New York, NY, USA, 882–905.

[125] Scott Wu. 2024. Meet Devin: The World's First AI Software Engineer. https://www.cognition-labs.com/introducing-devin.

[126] XAI. 2024. Open Release of Grok-1. https://x.ai/blog/grok-os.

[127] Boming Xia, Tingting Bi, Zhenchang Xing, Qinghua Lu, and Liming Zhu. 2023. An Empirical Study on Software Bill of Materials: Where We Stand and the Road Ahead. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2630–2642.

[128] Zhou Yang, Zhipeng Zhao, Chenyu Wang, Jieke Shi, Dongsun Kim, Donggyun Han, and David Lo. 2024. Unveiling Memorization in Code Models. In *2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE)*. IEEE

Computer Society, 856–856.

[129] Deborah Yao. 2023. One Year On, GitHub Copilot Adoption Soars. https://aibusiness.com/companies/one-year-on-github-copilot-adoption-soars. *AI Business* (27 6 2023). Accessed: Mar. 22 2024. https://aibusiness.com/companies/one-year-on-github-copilot-adoption-soars.

[130] Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2023. Large Language Model Unlearning. *CoRR* abs/2310.10683 (2023).

[131] Rui-Jie Yew. 2024. Break It 'Til You Make It An Exploration of the Ramifications of Copyright Liability Under a Pre-training Paradigm of AI Development. In *Proceedings of the Symposium on Computer Science and Law*. 64–72.

[132] Z. Yang *et al.*. 2023. Gotcha! This Model Uses My Code! Evaluating Membership Leakage Risks in Code Models. *arXiv preprint arXiv:2310.01166* (2023).

[133] Sheng Zhang and Hui Li. 2023. Code Membership Inference for Detecting Unauthorized Data Use in Code Pre-trained Language Models. *arXiv preprint arXiv:2312.07200* (2023).