

From Requirements to Code: Understanding Developer Practices in LLM-Assisted Software Engineering

Jonathan Ullrich, Matthias Koch

Fraunhofer IESE

Kaiserslautern, Germany

{jonathan.ullrich, matthias.koch}@iese.fraunhofer.de

Andreas Vogelsang

paluno – The Ruhr Institute for Software Technology

University of Duisburg-Essen

Essen, Germany

andreas.vogelsang@uni-due.de

Abstract—With the advent of generative LLMs and their advanced code generation capabilities, some people already envision the end of traditional software engineering, as LLMs may be able to produce high-quality code based solely on the requirements a domain expert feeds into the system. The feasibility of this vision can be assessed by understanding how developers currently incorporate requirements when using LLMs for code generation—a topic that remains largely unexplored. We interviewed 18 practitioners from 14 companies to understand how they (re)use information from requirements and other design artifacts to feed LLMs when generating code. Based on our findings, we propose a theory that explains the processes developers employ and the artifacts they rely on. Our theory suggests that requirements, as typically documented, are too abstract for direct input into LLMs. Instead, they must first be manually decomposed into programming tasks, which are then enriched with design decisions and architectural constraints before being used in prompts. Our study highlights that fundamental RE work is still necessary when LLMs are used to generate code. Our theory is important for contextualizing scientific approaches to automating requirements-centric SE tasks.

Index Terms—requirements, code generation, interview study, LLM, GenAI

I. INTRODUCTION

Implementation is the phase of the software engineering (SE) process in which a software design is translated into executable code [1]. All preceding SE activities culminate in the implementation phase. Previously elicited user and system requirements are documented and modeled in a design to be used for implementation. Requirements and design represent an abstraction of the system to be realized. They can be explicitly documented or exist implicitly in the developer’s head [1]. Although requirements usually reside in the problem space, i.e., describing *what* needs to be built, the design constitutes a bridge between the requirements and the implementation and defines *how* to build the system. Nevertheless, the design does not prescribe the exact steps to realize the system in code. To do so, developers must translate a descriptive view of the system into requirements and design into actionable implementation steps. In this process, developers must consider and reflect on the requirements.

Large language models (LLMs) learn powerful representations of natural language. Since code resembles natural languages in many ways [2], researchers applied LLMs to learn representations of code [3]. Besides general-purpose models, such as GPT-4 or Llama, LLMs more extensively trained on code, such as Codex [4], CodeLlama [5], and DeepSeek-Coder [6], are proficient at many code-related tasks, such as clone and defect detection, code repair, code completion, and code summarization as evaluated on many benchmarks [3], [7]. Tools such as GitHub Copilot offer IDE integrations for these models to further increase the ease of use by providing the model with context from the code repository.

Adopting LLMs for SE depends on the compatibility of these tools with the SE process [8]. Good SE processes have requirements and design phases before the implementation phase [1]. Much of the research, however, has not considered the role of requirements and design in LLM-assisted implementation. Indeed, requirements engineering and software design are highly underrepresented in studies [9]. On the contrary, the studies that focus on requirements as a starting point for code generation [10], [11] do not consider actual real-world requirements (e.g., functional requirements or user stories) but rather programming tasks from coding challenges.

To address the gap in the literature regarding the use of requirements in LLM-assisted implementation in practice, we set the following research objective:

How do practitioners incorporate requirements and design information in LLM-assisted implementation?

In particular, we aim to find out whether practitioners use requirements as part of the prompt in the interaction with code models. Further, as requirements and implementation are connected through software design, we investigate which information on the design is relevant to be included in the prompt. Doing so, we contribute to a better understanding of the role of requirements engineering and design in LLM-assisted SE.

Due to the exploratory nature of this study, we employed qualitative analysis based on interviews. We conducted interviews with 18 practitioners from 14 small to large companies

from 12 domains to find out how requirements and design information are reflected in the interaction with LLMs.

From the interviews, we derive a theory that describes how practitioners get *from requirements to code* when using LLMs. In this paper, we make the following contributions:

- We present a theory that describes **processes and interaction patterns** developers currently follow in practice when generating code based on requirements.
- Our theory also describes the **content** developers incorporate in the prompts to generate useful code, i.e., code that can be integrated in an existing code base.
- We provide our interview guide along with the codebook in a **replication package**¹.

The remainder of this paper is structured as follows: In Section II, we discuss work related to neural code intelligence and interaction with code models. We highlight the gap in the literature that this work aims to fill. In Section III, we describe our research method, the subjects of our study, and the data analysis process. In the following Section IV, we present the findings of our study and show the resulting process and content model. We discuss the relevance and impact of our findings in Section V. We conclude this paper in Section VI and highlight future work.

II. RELATED WORK

Neural Code Intelligence. Neural code intelligence [3], [7], [12] refers to the application of deep learning techniques to train models for understanding and generating code. A recent survey [3] identifies three evolutionary stages of neural code intelligence: embeddings, pretrained models, and LLMs. While embeddings [13] and pretrained models [14], [15] require task-specific training and fine-tuning, LLMs can be adapted to various tasks through prompting [16]. Several LLMs have been specifically trained in code, including Codex [4], PolyCoder [17], CodeGen [18], StarCoder [19], CodeLlama [5], and DeepSeek-Coder [6]. These models show state-of-the-art performance in tasks such as clone detection, defect detection, code repair, summarization, and generation [3], [7], [20]. In the following, we use the term *code models* to broadly refer to neural code intelligence models leveraging LLMs.

Generating Code from Natural Language Input. Beyond code understanding tasks, code models are predominantly used for code generation. Natural Language to Code (NL2Code) involves translating natural language descriptions into executable code. Code models have been evaluated on NL2Code tasks using various benchmarks [3], [7], ranging from fundamental programming challenges [21] to more complex SE tasks [22]. These evaluations primarily assess the quality of generated code based on diverse natural language inputs. Niu et al. [23] compare code models across 13 SE tasks, using natural language descriptions from the Concode dataset [24] and assessing their generated Java code. In Concode, natural language descriptions are embedded in Java classes and resemble

method comments (e.g., “adds a scalar to this vector in place” or “increment this vector” [24]). Similarly, Wang et al. [25] examine NL2Code at the method level by pairing code snippets from the CodeSearchNet dataset [26] with corresponding docstrings (e.g., “extracts video ID from URL” [26]).

Requirements as Natural Language Input. Although the above-mentioned studies focus on NL2Code, they do not consider the specifics of SE processes in which code generation occurs. In SE, requirements are the primary means of describing a system in natural language [1]. For code models to be effectively integrated into SE workflows, they must align with established processes such as RE [8], [27]. Liu et al. [10] introduced the “Requirements Text-based Code Generation” (ReCa) dataset to evaluate code generation based on requirements. However, their approach treats any textual program description as a requirement, disregarding the standard definition of requirements as “a condition or capability needed by a user to solve a problem or achieve an objective” [28]. The ReCa dataset, sourced from programming contest platforms, contains descriptions of programming tasks rather than genuine requirements (e.g., “You are given an array consisting of n integers. Your task is to find the maximum length of an increasing sub-array of the given array [...] Input: [...] Output: [...]” [10]). Mu et al. [11] sought to improve LLM-based code generation through requirements clarification but also treated programming task descriptions as requirements (e.g., “write a function to sort a list of elements” [11]). Thus, existing NL2Code evaluations have not considered requirements as defined in SE and RE communities.

Developer Practices with Code Models. Given that natural language and requirements can be ambiguous and often require clarification [29], NL2Code (or Requirements-to-Code, Req2Code) entails an interactive process with code models. Research has increasingly focused on integrating code models into user-friendly coding assistants. GitHub Copilot, a widely used in-IDE code completion tool, has been evaluated for its impact on developer productivity [30] and usability [31], [32]. Vaithilingam et al. [31] assess whether GitHub Copilot helps programmers complete tasks, finding that while it does not reduce task completion time, developers appreciate it as a useful starting point. Barke et al. [32] analyze user interactions with GitHub Copilot and identify two primary engagement patterns: using it to *accelerate* workflows and to *explore* possible solutions. Jiang et al. [33] conducted a user study in which participants used GenLine, a tool similar to GitHub Copilot, to solve programming tasks. Both Barke et al. and Jiang et al. show that users refine their natural language queries to optimize model responses, effectively developing a “syntax” for interacting with the models [32], [33].

Research Gap. The aforementioned studies evaluate NL2Code and coding assistants primarily for solving programming tasks. To the best of our knowledge, no research has explored the role of requirements and design artifacts in LLM-assisted implementation in industrial practice. Aligning code models with prior SE phases is essential for integrating AI into software development processes effectively [8]. To address this

¹<https://zenodo.org/records/15005613>

gap, this paper aims to provide a comprehensive study of how practitioners incorporate *real-world* requirements and design decisions when interacting with code models.

III. STUDY DESIGN

A. Research Method

The idea for this paper started as a broad research interest: *With increased use of code models in software engineering, what is the role of requirements? Are requirements, as natural language descriptions of a system, used as input for code models in practice?* Due to the nascent topic, we set a broad research focus instead of limiting our scope to predefined research questions. We employ an interview study [34] to qualitatively explore this research focus. Iteratively, we conducted interviews, analyzed the data, learned from the informants' statements, and applied the insights in subsequent interviews. In this process, we refined our research focus to include not only the role of requirements but also design information in the context of LLM-assisted software engineering. In this process, we narrowed our focus to the following research objective: *How do practitioners incorporate requirements and design information in LLM-assisted implementation?*

B. Interviews and Study Participants

In our study, we conducted interviews with 18 practitioners from 14 companies across 12 different domains between November 2024 and February 2025. We applied convenience sampling using our network to acquire interview partners. Many of the participants in this study stated that their company is pursuing internal programs to evaluate the usage of code models in their software engineering processes. The interview partners were selected based on their experience with code models in the implementation phase. During the selection of interview partners, a diverse set of roles and domains was included. As the focal point of this study lies in using code models to generate code that meets requirements, developers were primarily chosen as interview partners, as they are the ones primarily realizing the requirements in code. Other roles, however, have also been considered to gather a broad view of the topic. In larger companies, we were able to gather a holistic view of LLM-assisted software engineering by talking to team leads and a product owner responsible for the topic (P05, P06, P07). We tried to balance breadth and depth, exploring insights from many different companies in different domains, but also collecting multiple impressions from the same company if doing so deemed to provide a deeper understanding of the usage of code models at a company (P02 and P03, P06 and P07, P13 and P14, P15 and P16 were from the same companies). We stopped acquiring new interviewees once our theory was validated in subsequent interviews and no new themes emerged. This has been the case after we conducted the interviews scheduled with P15, P16, P17, and P18. Table I shows the participants' roles, the size of the company they work for (following the European definition of SMEs and larger firms), the application domain of the developed software, the type of requirements used by

the participants (standardized specifications containing (non-) functional requirements, user stories, or unstructured issues) and the interface for interacting with code models (either using a chat, e.g., ChatGPT, or an IDE-integration, e.g., GitHub Copilot).

The interviews, ranging from 30 to 60 minutes, were conducted online and recorded using Microsoft Teams. The interviews were conducted by the first author following a semi-structured interview guideline. In the first part of the interview, the interviewees were asked to describe their level of experience, the software engineering process at their company, and the artifacts (i.e., documented requirements or design decisions) created in the phases before using code models for implementation. The first part of the interview concluded with the participants describing their usual implementation workflow when tasked to realize a requirement without using code models. The second part of the interview focused on the use of code models in the implementation, with a dedicated focus on aligning the generated code with requirements and design decisions. Participants were asked to outline the "trigger" to start using code models, the content of the used prompts, and the process by which they ensure that the generated code fulfills a requirement. All interviews were transcribed. All transcripts were translated from German to English using DeepL, except for two interviews held in English.

C. Data Analysis and Theory Construction

We inductively created a theory based on interviewees' statements by exhaustive in-vivo coding. This resulted in 179 quotes as in-vivo codes. Iteratively, the in-vivo codes were clustered, creating higher-level themes. Themes that could not be confirmed by repeatedly being present in different interviews were discarded. The initial coding was done by one of the authors using the tool MaxQDA and exported to Excel. The coding of in-vivo codes and themes was reviewed by the other two authors regarding their consistency. Collectively, the authors named the themes and discussed the resulting theory. As, in software engineering, the same phenomenon can oftentimes be described from a behavioral and structural viewpoint, we also constructed our theory concerning these two viewpoints, resulting in a process model and content model of LLM-assisted implementation based on requirements and design decisions. The replication package shows the traces from in-vivo codes to themes to derived entities in the process and content model.

D. Threats to Validity

We took steps to mitigate common threats to validity in qualitative studies [35].

Internal validity refers to the threat that causal relationships or explanations are falsely drawn from the data. This can occur if the researchers' assumptions or participants' biased responses influence the interpretation of results. To mitigate this threat and to collect all relevant data, we recorded the interviews and transcribed them word by word. We annotated the transcripts with pointers to the respective positions in

TABLE I
INTERVIEW PARTICIPANTS

Participant	Role	Company Size	Application Domain	Requirements	Used Interaction
P01	Developer	Large	Service Management	User Stories	Chat
P02	Developer	Small	Construction Planning	Issues	IDE-Integration
P03	Developer	Small	Construction Planning	Issues	IDE-Integration
P04	Software Architect	Small	Smart City	Stand. Spec. + User Stories	Chat, IDE-Integration
P05	Product Owner	Large	Tax and Legal	User Stories	Chat, IDE-Integration
P06	Team Lead	Large	Enterprise Resource Planning	User Stories	Chat
P07	Team Lead	Large	Enterprise Resource Planning	User Stories	Chat
P08	DevOps Engineer	Small	Automotive	User Stories	Chat
P09	Data Scientist	Small	Smart City	User Stories	Chat
P10	Developer	Medium	Health	Stand. Spec. + User Stories	Chat
P11	Requirements Engineer	Medium	Aerospace	Stand. Spec. + User Stories	Chat, IDE-Integration
P12	Developer	Medium	Service Management	User Stories	Chat
P13	Developer	Medium	E-Commerce	User Stories	Chat, IDE-Integration
P14	Team Lead	Medium	E-Commerce	User Stories	Chat, IDE-Integration
P15	Team Lead	Small	Virtual Reality	User Stories	Chat, IDE-Integration
P16	Team Lead	Small	Virtual Reality	User Stories	Chat, IDE-Integration
P17	Developer	Large	Automotive	Stand. Spec. + User Stories	Chat
P18	Developer	Medium	Banking	Issues	IDE-Integration

the recording to trace back to the original conversation. We referred to in-vivo codes to ensure that interpretations were grounded in the actual data. To minimize the possibility of misunderstandings between interviewees and the researchers, the study goal was explained to the participants before the interview. During the interviews, we asked participants to show us intermediate artifacts such as user stories, prompt templates, and prompts, if possible.

External validity concerns the generalizability of our findings beyond the specific context of this study. While qualitative research does not aim for statistical generalization, we addressed this threat by including participants from diverse roles, industries, and levels of experience. This variety increases the likelihood that the identified themes apply across a broader range of settings. We tried to minimize sampling bias by including practitioners from 12 different domains.

Construct validity refers to the extent to which the interview questions and coding accurately capture the phenomena under investigation. To ensure construct validity and to mitigate researcher bias in the interpretation of the interviews, the coding was cross-validated by all authors and discussed in case of disagreements. Steps taken to further improve the reliability of the data included a review of the interview guideline. We also discussed derived themes and models to ensure consistent interpretation of constructs.

IV. FINDINGS

In the following, we present the main findings from our study. Firstly, we describe the *process* of using code models for implementation based on requirements and design information (Fig. 1). Secondly, we underpin the usage of code models in LLM-assisted SE by outlining the *content* that is given to a code model during this process (Fig. 2). We highlight the empirical evidence (i.e., participants mentioning a finding) in parentheses after each finding.

A. Process Model

Practitioners engage in a process of decomposing requirements artifacts into programming tasks followed by different approaches to implement the programming tasks. We discuss Fig. 1 in two steps: deriving programming tasks from requirements artifacts (the upper part) and the LLM-assisted implementation of programming tasks (the lower box).

Deriving Programming Tasks From Requirements Artifacts. Participants in our study stated that requirements artifacts, as traditionally documented, are too abstract to be fed into code models. Several interviewees tried using requirements artifacts as input for code models but discarded the generated code as they could not use it. P10 described this as follows: “If I say GPT write a command interface for a Keithley [multimeter], then I have in fact taken a requirement from my catalog [...] and I’ve never actually had anything ready to use come out of it”. P04 shared a similar experience: “I’ve found that if you try to implement entire requirements on this large level, depending on how they are formulated, then you end up with something, but you have to adapt a lot and that’s not efficient. That’s why there’s another ‘breaking down’ process before the prompting starts”.

Practitioners unanimously stated they first derive smaller units from the requirements artifacts. We refer to these smaller units as *programming tasks*. Requirements artifacts oftentimes leave implementation details unspecified. Programming tasks include more specific details on how to realize a requirement in code. When deriving programming tasks from requirements artifacts, these details are added: “The rough draft of the program, or how I want to structure it, I’ll do that myself. I think about what I need for the remaining part and then I think about what logic to put in, which function I need to call, how to connect to the database, and so on. This structure, I do all that myself. I don’t let ChatGPT do that, most of the time nothing comes out of it” (P06). Almost all interviewees confirmed this notion of deriving programming tasks from the

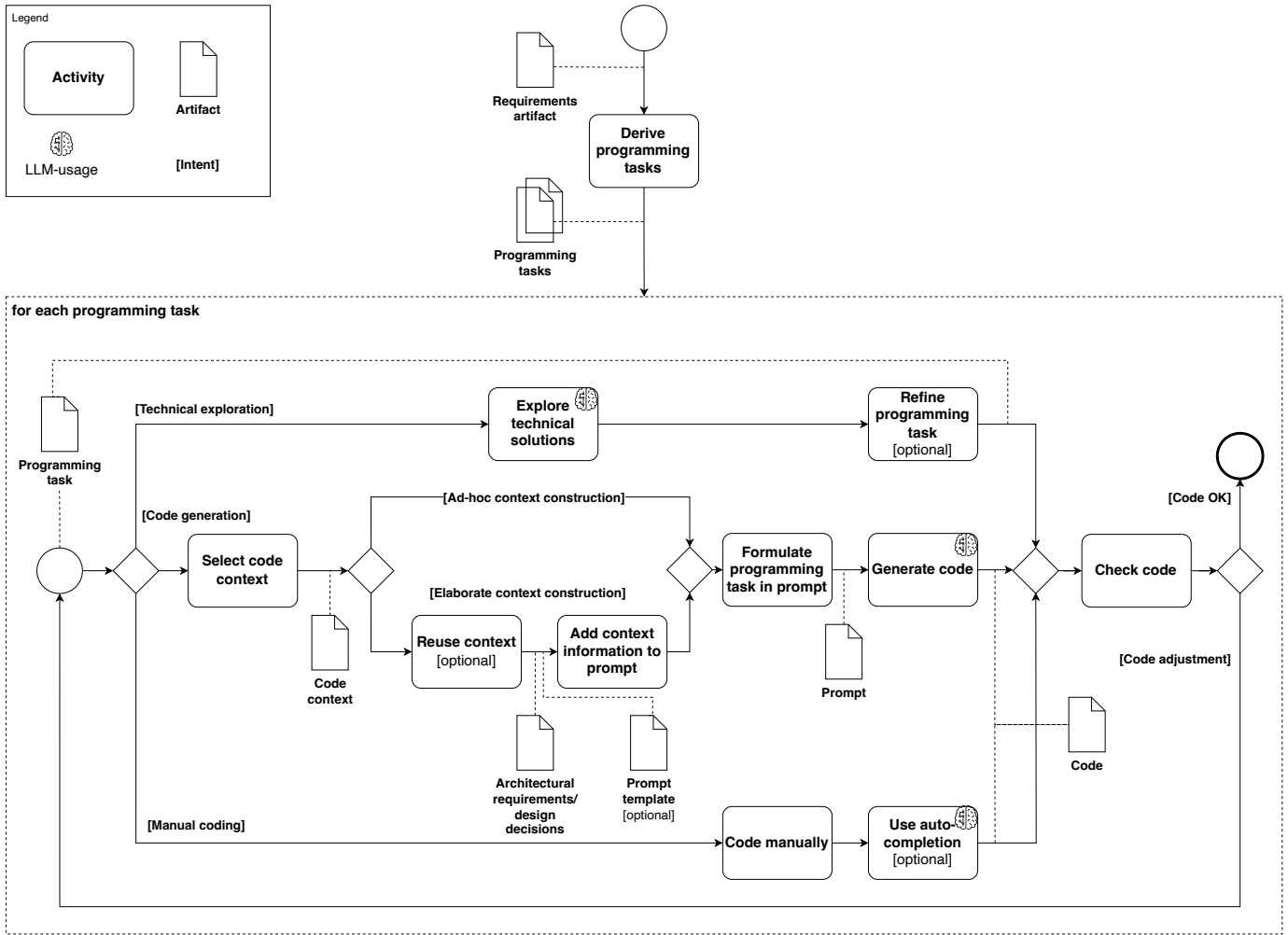


Fig. 1. Process model on LLM-assisted implementation

original requirements artifacts to use for code generation (P01, P04, P05, P06, P07, P08, P09, P10, P11, P13, P15, and P17).

P04, P05, and P06 argued that the process of refining requirements into programming tasks is not new to implementing with code models. P04 stated: “I think this is also the classic thought process of how to implement requirements before coding assistants existed. It’s not that easy to describe, because it happens subconsciously. You could now make a to-do list, and each of these to-dos would then become more or less a prompt with a feedback loop. That’s probably how it would be, but it’s not as if I would write all these steps down”.

P01, on the contrary, described a process of creating an explicit specification which defines the programming tasks for a user story: “The product owner creates the [user] stories [...] then, we go there later when we edit the story and write specs on how it should work technically”. The specification supports the construction of a prompt: “My specs say that if there is invalid input, I should handle it this way and this is one point in the spec that I’m trying to implement. Then I give the code model the existing code [and say] the behavior [from the spec] should be added” (P01).

P04, P08, and P13 suggested that the refinement of requirements could also be done by AI: “There would be a gap that could be closed by actually throwing the requirement to the model and then saying ‘first make a plan, write down these to-dos’ and then you can edit them again and then implement each one individually” (P04).

Takeaway 1: Practitioners do not use traditional requirements artifacts (such as user stories or functional requirements) as input for code models. Instead, they decompose and refine requirements artifacts into *programming tasks* to be used as input for LLM-assisted implementation.

Implementation of Programming Tasks. Our study reveals three different *modes* by which practitioners use code models to implement programming tasks: *technical exploration*, *code generation*, and *manual coding* (lower box in Fig. 1). Practitioners select and combine these modes depending on the familiarity of the developer with the programming task: “There is the side of implementation itself with things that you are

already familiar with. And there's the exploratory side, where you don't actually know what exactly you're supposed to do, and then I'll ask ChatGPT" (P10). Apart from technical know-how, experience using code models also determines the chosen mode, as illustrated by a statement made by P02: "I don't yet feel like I'm trying to do everything via code models. I still feel like a developer trying to do a lot myself and not thinking that code models could actually make some of my work easier. I'd have to practice that a bit more, I'd say".

The *technical exploration* mode is followed if a programming task is not clear enough to the developer or if the developer is uncertain about implementation details and context. LLMs support *technical exploration* by outlining solution approaches for a particular programming task. Our study participants use LLMs as an "intelligent search engine" (P15) or "Google replacement that gives you a quicker and more precise answer" (P12) during implementation. Based on the gained knowledge about possible solutions, the participants may further refine the programming task. Refinement of the programming task aims at improving the generated code: "The further down you go, the better the result ChatGPT gives you, i.e. the more precise you describe it, and the smaller the task is, the better it actually gets a suggestion that you can use as a guide and that is really good" (P07).

The *code generation* mode is followed when code models shall generate the majority of code for a certain programming task. It is initiated by selecting the code context relevant for code generation, i.e., selecting a section of existing code that needs to be considered when generating new code. Without the code context "there's too much knowledge on the outside that [the code model] needs", as P06 stated, "the simplest thing is to imagine which functions already exist where you don't have to develop them yourself and the model simply doesn't know that" (P06). Based on this initial context for the code model, practitioners pursue different steps depending on the task and the effort they are willing to invest in constructing more specific context. For simple tasks like generating 'boilerplate code' (P05), practitioners skip further context construction steps and craft the prompt based on the code context and the programming task (*ad-hoc context construction*, Fig. 1). P01 described her workflow in this regard as follows: "I often copy code that I already have, but now I would like [the model] to add this feature. So I already have working code, and now I would like that, when invalid input comes, it throws an error."

Contrary to this, there are cases where more context has to be given to the model to guide the output to generate useful code (*elaborate context construction*, Fig. 1). In these cases, information regarding design decisions and architectural constraints is added to the prompt. P05 exemplified this: "We have certain infrastructure requirements in our company that have to be met in a certain way, that of course the co-pilot or no model knows what specifications we have for how things are to be built". This context information has to be given in the prompt: "The context is super important. Sometimes you have to say [...] I have an app that is based on .NET or Java, and this is a web application, and this is the controller, and

the following error happens [...] anything that's not so obvious should be given, which could have an impact on the answer" (P11).

During the interviews, participants frequently referred to chat histories and prompt templates when describing their interaction with code models. P01 stated to fall back to previous chat histories when encountering a similar problem and wants to continue from there, reusing the previous context: "When I come across the same problem again, [...] I go back to the same prompt again and because I want it to remember all the stuff we've already discussed and now I want to add something". P06, P07, P08, and P11 use prompt templates for different implementation tasks. P13 and P17 set up 'Agents' or 'Custom GPTs' in their tooling, essentially pre-filling the context with information for repetitive tasks. Reusing chat histories, using prompt templates, and pre-filling context information reduce the effort to make the generated code useful in a certain context. Other interviewees confirmed this notion of different amounts of effort put into the formulation of the prompt (P06, P08, P09, P10, P13, P17).

All participants denied documenting the prompts for later reference. P11 replied: "I'm not doing that at the moment [...] It certainly wouldn't be bad, although I don't know whether you really look into it or not. So I think what helps more are prompt templates, i.e., saying that I have certain standard prompts and then maybe I can just fill in certain placeholders, like a form".

The *manual coding* mode is followed when developers prefer to start coding manually or the generated code needs to be adjusted. The majority of interviewees stated that the generated code always needs some manual adjustment. Some even said they prefer to start coding manually and then integrate generated code suggestions as they go (*Manual coding*, Fig. 1): "I either keep writing because the suggestion doesn't make sense or I just use the auto-complete suggestions and sometimes adjust them a bit" (P18). Regardless of the selected mode, the code must be inspected to check whether it realizes the programming task or the code needs to be adjusted. Practitioners iterate through the depicted process multiple times until they have implemented all programming tasks to realize a requirement.

Takeaway 2: Practitioners implement programming tasks by a combination of manual coding and code generation. When generating code, they try to reduce the effort of adding relevant context information to prompts by using prompt templates, reusing contexts from chat histories, or pre-filling context information.

Interaction Patterns. Practitioners reported different forms of iterating through the depicted process. The most commonly reported practice is using code models as a 'pair programmer' in a pattern we call *incremental code generation*. Practitioners follow all paths in the process shown in Fig. 1. They interleave technical exploration, manual coding, and code generation. In

this process, they incrementally construct context that they reuse when needed. This pattern is prevalent in reports of practitioners using a chat interface as it allows for iterative interaction and the reuse of the chat history (i.e., context information). Practitioners implement a programming task by manually adjusting the generated code or by reformulating the prompt.

In contrast to the aforementioned pattern, interviewees P02, P03, P05, and P18 described a pattern that can be summarized as *manual coding with intelligent auto-completion*. Here, code models have little impact on the practitioners’ workflow itself: Interviewees reported writing code as they normally would without code models but accept generated auto-completions and code suggestions, if applicable. Practitioners have mentioned this pattern mainly when using IDE-integrated code models. They spend most of the time in the ‘manual coding’ mode in Fig. 1. As the code model’s generated auto-completions are inserted where developers normally insert code, developers can assess the code with little effort. In this interaction pattern, developers rely less on the ability of code models to combine multiple implementation steps. Instead, they interact with code models at a lower level of abstraction when accepting code suggestions to implement a programming task.

Lastly, interviewees P04, P08, P13, and P14 stated that they use code generation whenever possible. We call this pattern *extensive code generation*, where practitioners offload a larger portion of the programming task to a code model. Practitioners who employ this form of interaction spend most of the time (re)formulating the prompt and checking whether the generated code fulfills the desired functionality: “I would say I don’t really write code myself anymore. I mostly iterate with AI” (P13).

Generally, our study indicates a trade-off between the effort invested in formulating prompts and the efficiency gained from generating code. P12 argued that: “If I choose every prompt super specifically, then I’ll get it right, but then I don’t need the AI either [...] If I had to think about every little thing myself, if I had to write it [in the prompt] myself, then I wouldn’t need the AI to help me in any way”.

Respondents using *manual coding with intelligent auto-completion* referred to lower level forms of interaction (closer to the code), whereas interaction with code models in *incremental* or *extensive code generation* might be more abstract, leveraging the capability of code models to combine multiple implementation steps in one output. From this finding, we derive the following takeaway:

Takeaway 3: The level of abstraction of a programming task varies depending on the interaction pattern. Practitioners break down requirements accordingly into programming tasks to accommodate a certain interaction pattern.

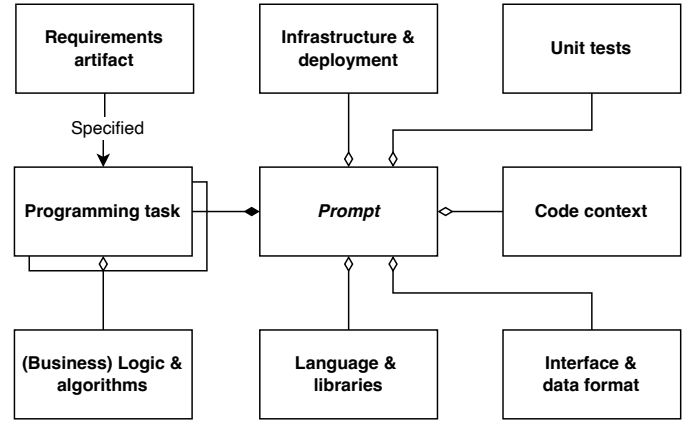


Fig. 2. Content model on LLM-assisted implementation

B. Content Model

In the content model (Fig. 2), we focus on the information from requirements and design decisions that is incorporated into the interaction with code models at different points in the process. One participant stated the motivation to add specific information about the constraints as follows: “You have to expand the context accordingly. All architectural boundary conditions really have to be documented in some way. And that goes right down to the skill sets of the individual developers, because they also influence architectural decisions [...] I think that would be the first step, i.e., to really expand the functional requirements with all the architectural conditions and then to expand them in the context of what is already there, so that it remains consistent with the previous code base” (P04). In this statement, P04 highlighted how (functional) requirements act as a starting point for LLM-assisted implementation but are only useful for code generation if enriched with information on specific design decisions. Other study participants also emphasized the importance of specifying architecturally significant requirements [36] and design decisions in the interaction with code models. Practitioners add information to the prompt context to ensure the generated code can be used, i.e., the code can be integrated into an existing code base. P06 stated this as follows: “All these boundary conditions that we have [...] flow into [the prompt]. For example, if it is already clear that we want to use a certain programming language or a framework, then that is included [...] because otherwise I might end up with a result that works on its own, but isn’t at all useful in our context”. P05 added to this: “When I say cloud-native, that’s a standing term for us, where we say we have a certain idea of what an application should look like, what environment it should run in, what basic systems we have, how things should fit together and how they have to be built. And these are just typical requirements that you can take into account, if you say I have a tool that generates code for me”.

Fig. 2 shows an overview of the content that interviewees repeatedly reported to include into the prompt for *code generation* (‘add context information to prompt’ in Fig. 1). From the interviews, we find that information related to *infrastructure*

and deployment (P04, P05, P06, P08), unit tests (P05, P06, P07, P10, P17), interface and data format (P03, P06, P09, P10), and language and libraries (P01, P02, P03, P04, P05, P07, P09, P12, P13, P15) is added to the prompt to ensure the generated code can be integrated into an existing code base. *Infrastructure and deployment* refers to “infrastructure requirements [...] that have to be met in a certain way that, of course, Copilot or no model knows” (P05). *Unit tests* and information about the *interface and data format* are added to the context as a way to guide the code models’ output. While the used *programming language and libraries* can oftentimes be inferred from the code context, study participants mentioned outdated information about libraries and missing information about domain-specific languages to limit the usefulness of the generated code. The *programming task* is the core component in the prompt and may include information about the *business logic and algorithms* to be used (as stated by P01, P02, P03, P09, P11, P13, and P16).

Takeaway 4: To craft a good prompt around a given programming task, practitioners add specific business logic or algorithmic information along with constraints on infrastructure and deployment, languages and libraries, interfaces and data formats, unit tests, and the code context in which the generated code shall be integrated.

V. DISCUSSION

Our study provides answers to our research objective: *How do practitioners incorporate requirements and design information in LLM-assisted implementation?* The findings suggest that constraints from architecturally significant requirements [36] or design decisions are directly used by developers in the interaction with code models; more abstract user-level requirements are not suited for direct input to generate useful code. Information from traditional requirements artifacts (such as user stories or standardized specifications containing functional requirements) is only indirectly incorporated in LLM-assisted implementation in practice. Practitioners manually decompose traditional requirements artifacts into concrete programming tasks to be included in the prompt. Design decisions and architecturally significant requirements (e.g., non-functional requirements) act as constraints for the generated code and are directly incorporated in the interaction with code models (i.e., specified in the prompt). There is, however, no consistent general approach to incorporating requirements and design information in LLM-assisted implementation. Our study contributes towards classifying different modes of interacting with code models. In the following, we discuss these findings in relation to existing evidence and how these impact research and practice.

Relation to existing evidence. Our findings can be related to several findings from previous studies. Barke et al. [32] identified two primary engagement patterns of developers when using LLMs for code generation: using LLMs to *accelerate* workflows and to *explore* possible solutions. Participants

in our study confirmed this notion. As shown in Figure 1, practitioners engage in activities to accelerate their work (‘generate code’ and ‘use auto-completion’) or to ‘explore technical solutions’. Our study puts Barke et al.’s findings into the perspective of the overall implementation process, including requirements and design artifacts.

Russo [8] emphasizes that “the adoption of Generative AI tools hinges largely on their successful integration with existing software development workflows” [8]. By outlining the process of LLM-assisted implementation, as reported by our study participants, we contribute to understanding when and how LLMs are used in existing SE processes in practice. Further, Russo found “concerns about the limitations of LLMs, such as their lack of knowledge about internal APIs and the need for human oversight, emphasize the importance of human expertise in utilizing these tools effectively” [8]. Our results confirm these findings in the sense that participants in our study also reported the need for searching and adding project-specific contextual information to prompts. They also unanimously reported that the resulting code must be checked and adapted in most cases. By outlining the content of a prompt in LLM-assisted implementation, we describe what information LLMs are missing that needs to be added manually.

Xu et al. [37] study the usage of IDE-integrated assistants for NL2Code tasks before the advent of LLMs. They evaluated an IDE plugin that takes a natural language query and outputs a list of generated code snippets and retrieved Stack Overflow answers. They already hypothesized that integrating context information into the queries is necessary for future work. Our results support this hypothesis and strengthen it further. Not only does the local code context play a crucial role, as hypothesized by Xu et al., but also additional information, such as architectural requirements and design decisions, is relevant for generating useful code.

We emphasize the role of programming tasks as technical specifications of how to implement a requirement in multiple steps in LLM-assisted implementation. A qualitative study by Liang et al. [38] focuses on “implementation design decisions”, i.e., decisions made by developers on how to implement a certain design when there are potential alternatives. Their study shows that developers constantly monitor high-level requirements and design decisions during implementation. However, they find that implementation decision-making processes are not consistent across developers. Our study outlines a process model and suggests interaction patterns used by developers. Similar to Liang et al. we also observed that practitioners do not follow a general structure to incorporate requirements-related information in their interaction with code models. Further research would need to identify conditions that influence developers to choose one or the other interaction pattern.

Impact for research. We believe our results contribute to contextualizing research on automatic requirements-centric SE tasks. By “contextualizing”, we mean evaluating approaches in light of current developer practices and the challenges they present. Our findings indicate that requirements, as doc-

umented in practice today, are rarely suitable for direct use as prompt inputs. Instead, practitioners decompose them into smaller tasks and enrich them with contextual information to generate useful code. Any approach that claims to automate a requirements-centric SE task (e.g., code generation [11], [39], test case generation [40], or model generation [41]) should explicitly discuss the representativeness of its requirements and how it addresses the specific demands for producing useful code in practice. Researchers should carefully distinguish between NL2Code and Req2Code (Requirements-to-Code) approaches.

Additionally, our findings suggest the need for further research on supporting practitioners in the manual steps outlined in our theoretical framework. One particularly interesting insight from our study is the significant role of decomposing requirements into programming tasks that are appropriately sized for LLM prompts. The key question is determining the optimal level of granularity for prompts. Research on this topic might benefit from comparing the level of granularity of the natural language input used in previous studies on NL2Code and the level of granularity of programming tasks derived from requirements in practice.

Another research direction emerging from our results involves helping practitioners construct and reuse context. LLM context, such as in chat histories, appears to be a crucial element to generating useful code. Practitioners recall similar tasks from the past, revisit the corresponding chat history, and reuse that context for new tasks. Systematically supporting the construction and reuse of prompt contexts (i.e., leveraging existing design documentation or stored chat histories for future tasks) remains an underexplored area of research.

A third research avenue is investigating the trade-off between investing effort in prompt engineering versus adapting the generated code. Participants in our study observed that the quality of generated code improves with more carefully crafted prompts. However, since the output is never perfect and always requires post-generation modifications, the effort spent on prompt engineering must be balanced against the effort required for code adaptation. One participant noted that when time is limited, they invest less effort in prompt engineering, which can lead to increased rework after code generation. Managing this trade-off remains an open question, closely related to determining what constitutes “good enough” requirements engineering [42]–[44].

Impact for practice. Based on our findings, we conclude that the vision of fully automated software engineering—where domain experts, without necessarily possessing technical expertise, simply articulate their needs and desires to an LLM, which then produces a complete application—is still a distant reality, particularly for more complex software. Our results indicate that while LLMs are actively used in practice for code generation, the process is preceded by several manual and creative steps that require expertise in requirements and software engineering.

Our interview participants confirmed that requirements, as documented today, are often too abstract and must be

broken down—an essential task in requirements engineering. Similarly, they emphasized that identifying relevant contextual information and incorporating it into a prompt is crucial for obtaining useful code. This ability, which is fundamental to requirements and software engineering, is difficult for non-technical experts to perform effectively.

Additionally, our findings reveal that prompts are currently not regarded as important engineering artifacts that require documentation, storage, traceability, or review. It would be valuable to assess the implications of this transient use of prompts. Since the quality and relevance of generated code depend on the input provided in the prompt, we hypothesize that understanding which prompt produced a given output and tracing the origins of prompt information will become increasingly relevant—especially in domains where accountability and correctness are critical.

Therefore, we conclude that requirements and software engineering activities and skills remain essential in contemporary software development projects, even when LLMs are used for code generation.

VI. CONCLUSION

In this paper, we propose a theory that contextualizes requirements, artifacts, and design decisions in LLM-assisted implementation. This theory is based on semi-structured interviews with 18 practitioners from 14 companies across 12 domains. We contribute to the understanding of LLM-assisted software engineering by describing how requirements and design decisions are incorporated into the LLM-assisted implementation process. Our findings show that practitioners break down requirements into programming tasks to serve as input for code models. They also integrate information from architectural requirements and design decisions to ensure the generated code aligns with the existing code base. Additionally, we outline different interaction patterns and the composition of prompts in LLM-assisted implementation. We discuss the implications of these findings for requirements engineering research and practice.

Our study provides a foundation for future research on key activities in LLM-assisted implementation. Further studies could explore the decomposition of requirement artifacts into programming tasks in greater detail or examine how context information can be retrieved from existing documentation based on the prompt content model proposed in this study.

ACKNOWLEDGMENTS

We thank our interview participants for their time and commitment to contributing to our study. We also thank the reviewers for their extensive feedback and suggestions to improve this paper.

REFERENCES

- [1] I. Sommerville, *Software Engineering*, 10th ed. Pearson, 2015.
- [2] A. Hindle, E. T. Barr, M. Gabel, Z. Su, and P. Devanbu, “On the naturalness of software,” *Communications of the ACM*, vol. 59, no. 5, pp. 122–131, 2016.

- [3] Q. Sun, Z. Chen, F. Xu, K. Cheng, C. Ma, Z. Yin, J. Wang, C. Han, R. Zhu, S. Yuan *et al.*, “A survey of neural code intelligence: Paradigms, advances and beyond,” *arXiv preprint arXiv:2403.14734*, 2024.
- [4] M. Chen, J. Twarek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman *et al.*, “Evaluating large language models trained on code,” *arXiv preprint arXiv:2107.03374*, 2021.
- [5] B. Roziere, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, R. Sauvestre, T. Remez *et al.*, “Code llama: Open foundation models for code,” *arXiv preprint arXiv:2308.12950*, 2023.
- [6] D. Guo, Q. Zhu, D. Yang, Z. Xie, K. Dong, W. Zhang, G. Chen, X. Bi, Y. Wu, Y. Li *et al.*, “DeepSeek-Coder: When the large language model meets programming—the rise of code intelligence,” *arXiv preprint arXiv:2401.14196*, 2024.
- [7] Y. Wan, Z. Bi, Y. He, J. Zhang, H. Zhang, Y. Sui, G. Xu, H. Jin, and P. Yu, “Deep learning for code intelligence: Survey, benchmark and toolkit,” *ACM Computing Surveys*, 2024.
- [8] D. Russo, “Navigating the complexity of generative AI adoption in software engineering,” *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 2024.
- [9] X. Hou, Y. Zhao, Y. Liu, Z. Yang, K. Wang, L. Li, X. Luo, D. Lo, J. Grundy, and H. Wang, “Large language models for software engineering: A systematic literature review,” *ACM Transactions on Software Engineering and Methodology*, vol. 33, no. 8, pp. 1–79, 2024.
- [10] H. Liu, M. Shen, J. Zhu, N. Niu, G. Li, and L. Zhang, “Deep learning based program generation from requirements text: Are we there yet?” *IEEE Transactions on Software Engineering*, vol. 48, no. 4, pp. 1268–1289, 2020.
- [11] F. Mu, L. Shi, S. Wang, Z. Yu, B. Zhang, C. Wang, S. Liu, and Q. Wang, “Clarifygpt: A framework for enhancing llm-based code generation via requirements clarification,” *Proceedings of the ACM on Software Engineering*, vol. 1, no. FSE, pp. 2332–2354, 2024.
- [12] M. R. I. Rabin, V. J. Hellendoorn, and M. A. Alipour, “Understanding neural code intelligence through program simplification,” in *29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*, 2021, pp. 441–452.
- [13] U. Alon, M. Zilberstein, O. Levy, and E. Yahav, “code2vec: learning distributed representations of code,” *Proceedings of the ACM on Programming Languages (POPL)*, vol. 3, pp. 1–29, 2019.
- [14] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang *et al.*, “Codebert: A pre-trained model for programming and natural languages,” *arXiv preprint arXiv:2002.08155*, 2020.
- [15] Y. Wang, W. Wang, S. Joty, and S. C. Hoi, “Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation,” *arXiv preprint arXiv:2109.00859*, 2021.
- [16] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [17] F. F. Xu, U. Alon, G. Neubig, and V. J. Hellendoorn, “A systematic evaluation of large language models of code,” in *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, 2022, pp. 1–10.
- [18] E. Nijkamp, B. Pang, H. Hayashi, L. Tu, H. Wang, Y. Zhou, S. Savarese, and C. Xiong, “Codegen: An open large language model for code with multi-turn program synthesis,” *arXiv preprint arXiv:2203.13474*, 2022.
- [19] R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chim *et al.*, “Starcoder: may the source be with you!” *arXiv preprint arXiv:2305.06161*, 2023.
- [20] M.-F. Wong, S. Guo, C.-N. Hang, S.-W. Ho, and C.-W. Tan, “Natural language generation and understanding of big code for AI-assisted programming: A review,” *Entropy*, vol. 25, no. 6, p. 888, 2023.
- [21] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le *et al.*, “Program synthesis with large language models,” *arXiv preprint arXiv:2108.07732*, 2021.
- [22] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. Narasimhan, “Swe-bench: Can language models resolve real-world github issues?” *arXiv preprint arXiv:2310.06770*, 2023.
- [23] C. Niu, C. Li, V. Ng, D. Chen, J. Ge, and B. Luo, “An empirical comparison of pre-trained models of source code,” in *IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 2136–2148.
- [24] S. Iyer, I. Konstas, A. Cheung, and L. Zettlemoyer, “Mapping language to code in programmatic context,” *arXiv preprint arXiv:1808.09588*, 2018.
- [25] S. Wang, M. Geng, B. Lin, Z. Sun, M. Wen, Y. Liu, L. Li, T. F. Bissyandé, and X. Mao, “Natural language to code: How far are we?” in *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2023, pp. 375–387.
- [26] H. Husain, H.-H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, “Codesearchnet challenge: Evaluating the state of semantic code search,” *arXiv preprint arXiv:1909.09436*, 2019.
- [27] A. Vogelsang, “From specifications to prompts: On the future of generative large language models in requirements engineering,” *IEEE Software*, vol. 41, no. 5, pp. 9–13, 2024.
- [28] “IEEE standard glossary of software engineering terminology,” *IEEE Std 610.12-1990*, pp. 1–84, 1990.
- [29] A. Vogelsang, A. Korn, G. Broccia, A. Ferrari, J. Fischbach, and C. Arora, “On the impact of requirements smells in prompts: The case of automated traceability,” in *47th IEEE/ACM International Conference on Software Engineering (ICSE-NIER)*, 2025.
- [30] A. Ziegler, E. Kalliamvakou, X. A. Li, A. Rice, D. Rifkin, S. Simister, G. Sittampalam, and E. Aftandilian, “Productivity assessment of neural code completion,” in *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, 2022, pp. 21–29.
- [31] P. Vaithilingam, T. Zhang, and E. L. Glassman, “Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models,” in *Chi conference on human factors in computing systems extended abstracts*, 2022, pp. 1–7.
- [32] S. Barke, M. B. James, and N. Polikarpova, “Grounded copilot: How programmers interact with code-generating models,” *Proceedings of the ACM on Programming Languages (OOPSLA)*, vol. 7, pp. 85–111, 2023.
- [33] E. Jiang, E. Toh, A. Molina, K. Olson, C. Kayacik, A. Donsbach, C. J. Cai, and M. Terry, “Discovering the syntax and strategies of natural language programming with generative language models,” in *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, 2022, pp. 1–19.
- [34] H. J. Rubin and I. S. Rubin, *Qualitative interviewing: The art of hearing data*. sage, 2011.
- [35] J. W. Creswell and J. D. Creswell, *Research design: Qualitative, quantitative, and mixed methods approaches*. Sage publications, 2017.
- [36] L. Chen, M. Ali Babar, and B. Nuseibeh, “Characterizing architecturally significant requirements,” *IEEE Software*, vol. 30, no. 2, pp. 38–45, 2013.
- [37] F. F. Xu, B. Vasilescu, and G. Neubig, “In-IDE code generation from natural language: Promise and challenges,” *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 31, no. 2, pp. 1–47, 2022.
- [38] J. T. Liang, M. Arab, M. Ko, A. J. Ko, and T. D. LaToza, “A qualitative study on the implementation design decisions of developers,” in *IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 435–447.
- [39] B. Wei, “Requirements are all you need: From requirements to code with LLMs,” in *IEEE 32nd International Requirements Engineering Conference (RE)*. IEEE, 2024, pp. 416–422.
- [40] C. Arora, T. Herda, and V. Homm, “Generating test scenarios from NL requirements using retrieval-augmented LLMs: An industrial study,” in *IEEE 32nd International Requirements Engineering Conference (RE)*, 2024, pp. 240–251.
- [41] A. Ferrari, S. Abualhaija, and C. Arora, “Model generation with LLMs: From requirements to UML sequence diagrams,” in *IEEE 32nd International Requirements Engineering Conference Workshops (REW)*, vol. 24. IEEE, 2024, p. 291–300.
- [42] S. Fricker, T. Gorschek, C. Byman, and A. Schmidle, “Handshaking: Negotiate to provoke the right understanding of requirements,” *IEEE Software*, 2019.
- [43] J. Frattini, J. Fischbach, D. Fucci, M. Unterkalmsteiner, and D. Mendez, “Measuring the fitness-for-purpose of requirements: An initial model of activities and attributes,” in *32nd International Requirements Engineering Conference (RE)*. IEEE, 2024, pp. 398–406.
- [44] H. Femmer and A. Vogelsang, “Requirements quality is quality in use,” *IEEE Software*, vol. 36, no. 3, p. 83–91, May 2019. [Online]. Available: <http://dx.doi.org/10.1109/MS.2018.110161823>