# A Short Survey on Formalising Software Requirements with Large Language Models[*]

Arshad Beg[1][0009−0004−6939−0411], Diarmuid O'Donoghue[1][0000−0002−3680−4217], and Rosemary Monahan[1][0000−0003−3886−4675]

Department of Computer Science, Maynooth University, Ireland
`arshad.beg@mu.ie, diarmuid.odonoghue@mu.ie, rosemary.monahan@mu.ie`

**Abstract.** This paper presents a focused literature survey on the use of large language models (LLM) to assist in writing formal specifications for software. A summary of thirty-five key papers is presented, including examples for specifying programs written in Dafny, C and Java. This paper arose from the project "VERIFAI: Traceability and verification of natural language requirements" that addresses the challenges in writing formal specifications from requirements that are expressed in natural language. Our methodology employed multiple academic databases to identify relevant research. The AI-assisted tool Elicit facilitated the initial paper selection, which were manually screened for final selection. The survey provides valuable insights and future directions for utilising LLMs while formalising software requirements.

## 1 Introduction

Relying on natural language in software requirements often leads to ambiguity. In addition, requirements which are not expressed in a formal mathematical notation cannot be guaranteed through formal verification techniques as required to meet standards e.g. [6,46,4] in safety critical software. Expressing requirements in formal notation requires training in the domain of requirements engineering, as well as knowledge of formal notation and associated proof methods, often increasing the software development cycle time by a factor of 30% [24]. In this project, we aim to ease the burden of writing specification, helping to bridge the gap between the need for formal verification techniques and their lack of use in the software industry (due to the fast pace of the industry environment).

Formalising software requirements ensures clarity, correctness and verifiability, and requires formal specification languages, logic and verification techniques, such as theorem proving and model-checking, to guarantee the correctness of the software system under construction. Like all other fields, the development of Large Language Models (LLMs) has opened a world of opportunities where we can exploit their power to generate formal requirements and accompanying specifications.

In this paper, we present the results of our structured literature review which examines how large language models are currently used to assist in writing formal specifications. Main research questions for conducting systematic literature review on the topic are as follows:

**RQ1:** What methodologies leverage Large Language Models (LLMs) to transform natural language software requirements into formal notations?

**RQ2:** What are the emerging trends and future research directions in using LLMs for software requirements formalisation?

---

We represent a motivating example of generating formal requirements and specifications in section 2. Section 3 presents our methodology for the literature review. A summary of contribution from each paper addressing **RQ1** is provided in section 4. Section 5 presents future directions based on our findings and the discussion on chain of thought and prompt engineering in sub-section 5.1, addressing **RQ2**. Concluding remarks are provided in Section 6.

## 2      Example of Generating Assertions

Dafny is a verification-aware programming language that has native support for expressing formal specifications and is equipped with a static program verifier to automatically verify implementations against specifications using deductive verification. Mugnier et al. [37] present a framework named Laurel to generate Dafny assertions using LLMs, using two domain-specific prompting techniques. The first technique locates the position in code where an assertion, which provides part of the formal specifications, is missing. This is done through analysis of the Dafny verifier's error message. A placeholder is inserted at the particular location where the assertion is missing. The second technique involves the provision of example assertions from a codebase. Here, Laurel is able to generate over 50% of the required helper assertions, making it a viable approach to deploy, for automating the program specification and verification process.

### 2.1      Generating Dafny Assertions

We illustrate an example using a Dafny lemma as reported in [37], where a helper assertion is needed in the verification process. Lemmas are used in Dafny to specify properties that may be used in the verification process. The lemma, in this example (the full body of Lemma is included in Appendix 8.1, ensures that the integer and fractional parts are correctly extracted while parsing a decimal string.

Explaining the lemma presented, we have:

**Precondition** (**requires** clause): The function assumes that every character in `s1` is a digit ('0' to '9').

**Postcondition** (**ensures** clause): The function guarantees that when `ParseDecStr` is applied to the concatenation of `s1` and `s2`, the first part of the result remains `s1`, and the second part contains `"." + s2`.

**Base Case**: If `s1` contains a single character, the function asserts that parsing `"."+s2` leads to an empty integer part.

**Recursive Case**: The function calls itself with the tail of `s1` to process it recursively. However, to help the verifier understand the transformation, an assertion is added:

```
assert s1 + "." + s2 == [s1[0]] + (s1[1..] + "." + s2);
```

This assertion explicitly states how the string is decomposed and aids the SMT solver in proving correctness.

**Role of the Helper Assertion:**

Without the assertion, the Dafny verifier struggles to establish the correctness of the postcondition due to the complexity of reasoning about string concatenation. The assertion serves as an intermediate step, breaking down the transformation into a form that is easier for the solver to handle.

This example motivates our work by demonstrating the importance of inserting specification in the form of helper assertions in Dafny proofs. Tools like Laurel [37] aim to automate this process by leveraging Large Language Models to determine the relevant assertions.

## 3   Methodology for Literature Review

To conduct a structured and thorough review of literature on Natural Language Processing (NLP), Large Language Models (LLMs), and their use in software requirements, the following approach is followed. Several academic databases, including IEEE Xplore, ACM Digital Library, Scopus, Springer Link, and Google Scholar, are searched using specific keywords. The core search terms include "NLP," "LLMs," and "Software Requirements," with broader terms such as "specification," "logic," "verification," "model checking," and "theorem proving" used to expand the scope. The number of results differs notably across databases. For example, IEEE Xplore returns 17 peer-reviewed articles, Scopus lists 20, Springer Link filters to 595, ACM Digital Library provides 1,368 results, and Google Scholar shows 14,800 references since 2021. These discrepancies highlight the importance of applying precise selection methods to extract the most relevant studies.

To streamline the process of locating strong contributions, the AI-powered tool Elicit [11] is used. Elicit supports the literature review by offering summarised content and DOIs for suggested papers. While it helps reduce the manual workload during the initial phase, every suggested paper in this review is manually reviewed to ensure its relevance. This ensures that the final list excludes any unrelated or off-topic material. After the initial filtering, a manual review is performed to confirm the relevance and quality of each paper. Abstracts are first assessed to judge suitability. If the abstract lacks clarity or depth, a further examination of the full text is conducted.

Abstracts are closely read, and when necessary, the full text is reviewed using the following exclusion and inclusion criteria:

**Inclusion Criteria:** Studies are included if they offer meaningful theoretical or empirical insights related to NLP, LLMs, and their application in software requirements. This includes topics like specification, formal logic, verification, and formal methods.

**Exclusion Criteria:** Papers are excluded if they show in-sufficient relevance to the intersection of NLP/LLMs and software requirements, or if their abstracts or full texts lack sufficient detail. Non-peer-reviewed materials, duplicates, and items suggested by Elicit but deemed irrelevant after manual review are also removed.

## 4   Formalising Requirements through LLMs

This section provides an overview of the literature surveyed. A summary table is included in Appendix 8.2, listing each citation along with the tool, technique, or framework developed or analysed in the paper, accompanied by a brief description of its contribution. Figure 1 presents a tree-structured overview of the studied literature.

Table 1 presents a classification of the surveyed literature based on their methodological approach. Most works fall under the "Prompt-only" category, where LLMs are used without further tuning. Some studies involve human-in-the-loop or iterative prompting, while others fine-tune LLMs for improved performance. A set of works integrates LLMs with verifiers, and a few adopt neuro-symbolic methods combining LLMs with SMT solvers or theorem provers. Other categories include baseline/manual approaches using controlled natural language, meta-analyses and tool support
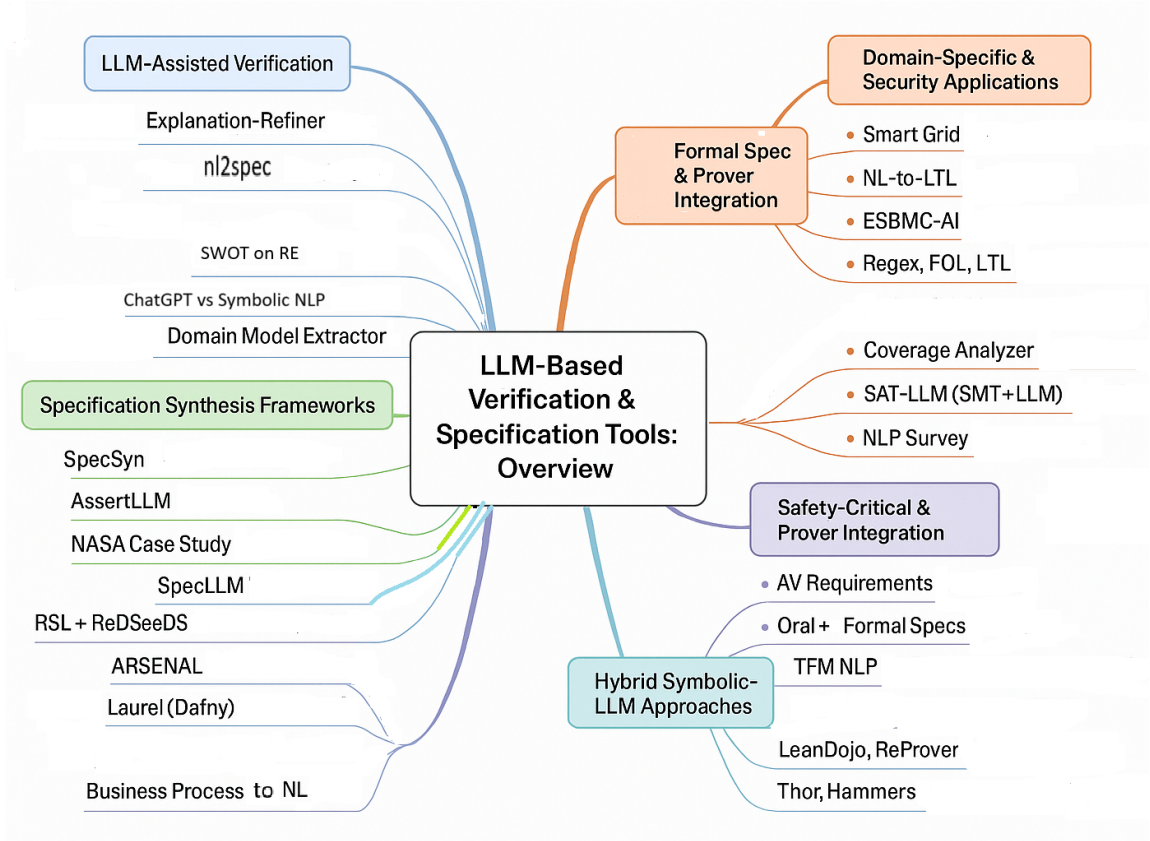
Fig. 1: LLM-Based Verification and Specification Literature Overview

studies, and a single work proposing IDE integration. Table 1 is a condensed version of Table 5 found in the Appendix. To provide readers with a more comprehensive reference, we have also included Table 4 in the Appendix. This table integrates the developed tools, their methodological classifications, and a brief summary of each work, offering a clear overview of the purpose, approach, and contributions of each paper at a glance.

## 4.1   LLM-Assisted Verification and Specification Tools

The paper [10] proposes using LLMs, like GPT-3.5, to verify code by analysing requirements and explaining whether they are met. The work [48] provides verification and refinement of natural language explanations by making LLMs and theorem provers work together. A neuro-symbolic framework i.e. Explanation-Refiner is represented. LLMs and theorem provers are integrated together to formalise explanatory sentences. The theorem prover then provides the guarantee of validated sentence explanations. Theorem prover also provides feedback for further improvements in NLI (Natural Language Inference) model. Error correction mechanisms can also be deployed by using the tool Explanation-Refiner. Consequently, it automatically enhances the quality of expla-

Table 1: Classification of Surveyed Literature by Methodology

| References | Classification |
|---|---|
| [6], [46], [10], [48], [2], [49], [56], [18], [52], [7], [54], [55], [50], [19], [61], [35], [10], [48], [59], [61], [52] | Prompt-only |
| [48], [10], [56], [50], [7] | Prompt + Iterative / Human-in-loop / CoT |
| [4], [44], [57], [27], [26], [56], [57] | Fine-tuned |
| [9], [35], [17], [51], [54], [61], [7] | Verifier-in-loop |
| [16], [39], [48] | Neuro-symbolic (LLM + SMT/Theorem Prover) |
| [31], [45], [44], [24], [31], [45] | Baseline / Manual / Controlled NL |
| [27], [3], [12], [34], [3] | Meta-analysis / Tool Support |
| [28] | IDE Integration Proposal |

nations of variable complexity. The work [13] evaluated GPT-4o's ability to generate specifications for C programs that can be verified using VeriFast, a static verifier based on separation logic. Their experiments, which use different user inputs and prompting techniques, show that while GPT-4o's specifications maintain functional behaviour, they often fail verification and include redundancies when verifiable.

## 4.2   Frameworks for Specification Synthesis

The work [9] details about nl2spec, a framework that leverages LLMs to generate formal specifications from natural language, addressing the challenge of ambiguity in system requirements. Users can iteratively refine translations, making formalization easier. The work [9] provides an open-source implementation with a web-based interface. An automatic synthesis of software specifications is provided through LLMs in [35]. Work [35] proposed SpecSyn framework that uses an advanced language model to automatically generate software specifications from natural language text. It treats specification generation as a sequence-to-sequence learning task and outperforms previous tools by 21% in accuracy, extracting specifications from both single and multiple sentences.

The paper [39] introduced Req2Spec, an NLP-based tool that analyses natural language requirements to create formal specifications for `HANFOR`, a large-scale requirements and test generation tool. Tested on 222 automotive software requirements at BOSCH, it correctly formalized 71% of them. The work [34] represents a novel framework named SpecGen to generate specifications through LLMs. Two phases are applied. First phase is about having prompts in conversational style. Second phase is deployed where correct specifications are not generated. Here, four mutation operators are applied to ensure the correctness of the generated specifications. Two benchmarks i.e. SV-COMP and SpecGen are used. Verifiable specifications are generated successfully for 279 out of 384 programs, making [34] a viable approach.

## 4.3   Domain-Specific and Security-Focused Applications

AssertLLM tool is presented in [14]. The tool generates assertions to do hardware verification from design specifications, exploiting three customised LLMs. It is done in three phases, first un-

derstanding specifications, mapping signal definitions and generating assertions. The results show that AssertLLM produced 89% correct assertions with accurate syntax and function. The work [16] reports on formal verification of NASA's Node Control Software natural language specifications. The software is deployed at International Space Station. Errors found in the natural language requirements are reported by the authors with a commentary on lessons learnt.

SpecLLM [31] explores the space of generating and reviewing VLSI design specifications with LLMs. The utility of LLMs is explored with the two stages i.e. (1) **generation** of architecture specifications from scratch and from register transfer logic (RTL) code; and (2) **reviewing** these generated specifications. In [49], the potential and power of LLMs is exploited for smart grid requirement specifications improvement. Here, the performance of GPT-4o and Claude 3.5 Sonnet is analysed through f1-scores, achieving in range of 79% - 94%.

### 4.4   Formal Specification Translation and Evaluation

In [27], symbolic NLP and ChatGPT performance is compared while generating correct formal contracts written in the Java Modelling Language (JML) which have been extracted from natural language specifications. The paper [56] reports the translation between NL and Linear Temporal Logic (LTL) formulas through the use of LLMs. Dynamic prompt generation and human interaction with LLMs are amalgamated to deal with the mentioned challenges. Unstructured natural language requirements are converted to NL-LTL pairs. The approach achieved up to 94.4% accuracy on publicly available datasets with 36 and 255,000 NL-LTL pairs. The primitive work [42] described automatic translation from natural language sentences to temporal logic, in order to deploy formal verification of the requirements.

### 4.5   Safety-Critical and Prover-Integrated Applications

The work [58] introduced LeanDojo, an open-source toolkit that enables programmatic interaction with Lean theorem prover. Using LeanDojo's extracted data, [58] developed ReProver, a retrieval-augmented LLM-based prover. Thor [25] is a framework developed to integrate language models with theorem provers. It improved the accuracy from 39% to 57% using the PISA dataset and also outperformed previous works on the MiniF2F dataset. The paper [19] explores integrating Co-pilot with formal methods. The integration of Copilot and formal methods is proposed through development of IDE containing language servers. Granberry et al. [18] explored combining LLMs with symbolic analysis to generate specifications for C programs. They enhanced LLM prompts using outputs from PathCrawler and EVA to produce ACSL annotations.

### 4.6   Hybrid Symbolic-LLM Methodologies and Future Outlook

SAT-LLM, a unique framework to remove conflicting requirements is represented in [15]. It integrated Satisfiability Modulo Theories (SMT) solvers with LLMs. SAT-LLM performed better than ChatGPT alone, identifying 80% of conflicts with a Precision of 1.00, Recall of 0.83, and an F1 score of 0.91. [2] outlines key research directions for the stages of software requirement engineering, conducts a SWOT analysis, and share findings from an initial evaluation.

The purpose of Dafny is to automate proofs by outsourcing them to an SMT solver. [37] presented a framework named Laurel to generate Dafny assertions using LLMs. Laurel was able to generate over 50% of the required helper assertions. [12] used input of error messages, variable names, procedure documentation and user questions. They discussed available literature for generating assertions by synthesising sentences in testing phase.

### 4.7 Historical and Foundational Work

The paper [44] introduced a model-based language (Requirements Specification Language - RSL). The framework functionality is integrated as development platform named ReDSeeDS. A similar work is reported in [17] which describes the ARSENAL framework and methodology designed to perform automatic requirements specification extraction from natural language. An interesting work is presented in [28]. It includes generation of natural language from business process models. The generated natural language is found complete and more understandable. In primitive work of 1996 [45], software requirements were expressed in a limited set of natural language referred to as controlled natural language. The primitive work in the domain is about RML [20]. RML bundled with features of writing requirements, which are based on conceptual model.

### 4.8 Evaluation of the literature surveyed

The literature surveyed in section 4 suggests that assertion generation currently shows higher reliability than full contract synthesis. Tools like Laurel and AssertLLM demonstrate success rates above 50% and 89%, respectively, when generating helper assertions or verifying hardware designs. These assertions typically focus on specific program points or behaviors, making them easier for LLMs to manage. In contrast, full contract generation, such as translating natural language to formal Java Modeling Language (JML) contracts or complete temporal logic formulas, often results in partially correct or unverifiable outputs. For example, efforts using GPT-4o to produce VeriFast-compatible specifications found that generated contracts often failed formal verification checks despite being functionally reasonable.

There is a clear trend where tasks involving smaller, well-scoped units like assertions yield more accurate results from LLMs. This is likely due to the limited context and reduced complexity compared to full contract generation. Full formal specifications, especially across multiple sentences or involving system-wide properties, require greater abstraction and contextual understanding. While frameworks like SpecGen and SpecSyn have made progress, their outputs often need iterative refinement or mutation operators to reach acceptable accuracy. This contrasts with tools like AssertLLM or Laurel, where results are immediately usable or require minimal correction.

For future outlook, this short paragraph is based on our discussion of Section 5 of the paper as well. We can say that the combination of LLMs with formal methods (e.g., theorem provers or SMT solvers) shows promise in boosting the reliability of both assertion and contract synthesis. Neuro-symbolic frameworks like Explanation-Refiner or SAT-LLM provide structured error correction and logic validation, improving overall quality. Similarly, iterative approaches like nl2spec and prompt engineering techniques (e.g., Chain-of-Thought prompting) help refine complex contract outputs. Overall, while assertions are currently more robustly supported, the research trajectory points toward improving full contract reliability through hybrid systems, fine-tuned prompts, and verification-aware model integration.

## 5 Future Directions

In this section, we outline prospective directions informed by the literature review. Much of the literature in this review employed queries containing a problem description and some instructions to achieve a desired outcome. Such querying of LLMs without training or examples of the current task is typically referred as zero-shot prompting and shows excellent performance on many tasks

[26]. Surprisingly, they also showed that the performance of LLM on some challenging problems can be improved by encouraging the LLM to reason using intermediate steps through a simple addition to problem prompts ("Lets think step by step").

### 5.1   Advanced Prompt Engineering

Beyond this approach is one-shot prompting that includes an example of a solved problem to guide the LLM into generating the desired output [32]. This can be extended to few-shot prompting where a number of differing examples guide the LLM. But improved results are not assured as some studies e.g. [60] show that zero-shot can outperform the few-shot case [61]. [22] reviewed the evolution of prompt engineering in LLMs, including discussions on self-consistency and multimodal prompt learning. It also reviewed the literature related to adversarial attacks and evaluation strategies for ensuring robust AI interactions.

Chain of Thought (CoT) [52] prompting involves a sequence of prompts producing intermediate results that are generated by the LLM and used to drive subsequent prompting interactions. These orchestrated interactions can improve LLM performance on tasks requiring logic, calculation and decision-making in areas like math, common sense reasoning, and symbolic manipulation. CoT requires the LLM to articulate the distinct steps of its reasoning, by subdividing larger tasks into multi-step reasoning stages, acting as a precursor for subsequent stages. But the CoT approach may require careful analysis when used with larger LLM offering long input contexts. This is because of the lost-in-the-middle problem where LLM show a U-shaped attention bias [22] and can fail to attend to information in the middle of the context window.

PromptCoT [59] enhanced the quality of solutions for diffusion-based generative models by employing the CoT approach. The computational cost is minimised through adapter-based fine-tuning. Prompt design is explored in detail in [1]. It discussed Chain-of-Thought and Reflection techniques, along with best practices for structuring prompts and building LLM-based agents.

Besta et al. [5] introduced the concept of reasoning topologies, examining how structures such as Chains, Trees, and Graphs of Thought improve LLM reasoning. They also proposed a taxonomy of structured reasoning techniques, highlighting their influence on both performance and efficiency. Structured Chain-of-Thought (SCoT) prompting was proposed by [30] to enhance code generation by incorporating structured programming principles. This approach significantly improved the accuracy and robustness of LLM-based code synthesis compared to standard CoT methods. Building on the theme of automation, [50] introduced Automate-CoT, a technique for automatically generating and selecting rational chains for CoT prompting. By minimising dependence on human annotations, it enabled more flexible adaptation of CoT strategies across diverse reasoning tasks. Complementing these efforts, [53] presented a prompt pattern catalog that offered reusable design patterns to optimise LLM interactions, thereby refining prompt engineering practices for a wide range of applications. Additionally, [55] proposed Reprompting (Gibbs sampling-based algorithm) for discovering optimal CoT prompts. The proposed prompting technique consistently outperformed human-crafted alternatives and demonstrated high adaptability across various reasoning benchmarks.

Retrieval Augmented Generation (RAG) [29] supplements problem information with specifically retrieved information and is often used in knowledge intensive tasks. This helps ensure the LLM attends specifically to the retrieved information when addressing the users prompt. LLM model selection (chat vs reasoning) and fine tuning such as with LoRA [23] remain among a growing number of possibilities for exploration.

Based on the literature survey conducted, we sketch one line research agenda: in VERIFAI, we aim to improve the techniques that bridge the gap between informal natural language description

and rigorous formal specifications, through refinement of prompt engineering, the incorporation of chain-of-thought reasoning and the development of hybrid neuro-symbolic approaches.

## 6   Conclusions

The role of large language models in formalising software requirements is surveyed in this paper. The key contribution of the selected papers is on bridging the gap between informal natural language descriptions and rigorous formal specifications. We can enhance requirement formalisation through automating translations. The accuracy of translated ones is of key importance. We can deploy iterative refinements to improve the accuracy and correctness of the generated requirements. While LLMs is contributing significantly in improving development cycle, the challenges of ambiguity resolution, verification and domain adaptation shall be the focus areas of future research.

In future research, the refinement of prompt engineering techniques, the incorporation of chain-of-thought reasoning and the development of hybrid neuro-symbolic approaches are the key areas to look at. Our survey concludes with the remarks that the better collaboration with industry and academia will further enhance the domain.

## 7   Acknowledgements

## References

1. Amatriain, X.: Prompt design and engineering: Introduction and advanced methods (2024), https://arxiv.org/abs/2401.14423
2. Arora, C., Grundy, J., Abdelrazek, M.: Advancing requirements engineering through generative ai: Assessing the role of llms (2023), https://arxiv.org/abs/2310.13976
3. Arora, C., Sabetzadeh, M., Briand, L., Zimmer, F.: Extracting domain models from natural-language requirements: approach and industrial evaluation. In: Proceedings of the ACM/IEEE 19th International Conference on Model Driven Engineering Languages and Systems. p. 250–260. MODELS '16, Association for Computing Machinery, New York, NY, USA (2016). https://doi.org/10.1145/2976767.2976769, https://doi.org/10.1145/2976767.2976769
4. Bell, R.: Introduction to iec 61508. In: Proceedings of the 10th Australian Workshop on Safety Critical Systems and Software - Volume 55. p. 3–12. SCS '05, Australian Computer Society, Inc., AUS (2006)
5. Besta, M., Memedi, F., Zhang, Z., Gerstenberger, R., Blach, N., Nyczyk, P., Copik, M., Kwasniewski, G., Müller, J., Gianinazzi, L., Kubicek, A., Niewiadomski, H., Mutlu, O., Hoefler, T.: Topologies of reasoning: Demystifying chains, trees, and graphs of thoughts. CoRR **abs/2401.14295** (2024). https://doi.org/10.48550/arXiv.2401.14295, https://doi.org/10.48550/arXiv.2401.14295
6. Brosgol, B.: Do-178c: the next avionics safety standard. Ada Lett. **31**(3), 5–6 (Nov 2011). https://doi.org/10.1145/2070336.2070341, https://doi.org/10.1145/2070336.2070341
7. Caglayan, B., Wang, M., Kelleher, J.D., Fei, S., Tong, G., Ding, J., Zhang, P.: Bis: Nl2sql service evaluation benchmark for business intelligence scenarios. In: Service-Oriented Computing: 22nd International Conference, ICSOC 2024, Tunis, Tunisia, December 3–6, 2024, Proceedings, Part II. p. 357–372. Springer-Verlag, Berlin, Heidelberg (2024). https://doi.org/10.1007/978-981-96-0808-9_27

8. Casadio, M., Dinkar, T., Komendantskaya, E., Arnaboldi, L., Daggitt, M.L., Isac, O., Katz, G., Rieser, V., Lemon, O.: Nlp verification: Towards a general methodology for certifying robustness (2025), https://arxiv.org/abs/2403.10144

9. Cosler, M., Hahn, C., Mendoza, D., Schmitt, F., Trippel, C.: nl2spec: Interactively translating unstructured natural language to temporal logics with large language models (2023), https://arxiv.org/abs/2303.04864

10. Couder, J.O., Gomez, D., Ochoa, O.: Requirements verification through the analysis of source code by large language models. In: SoutheastCon 2024. pp. 75–80 (March 2024). https://doi.org/10.1109/SoutheastCon52093.2024.10500073

11. Elicit: Elicit - the ai research assistant. https://www.elicit.com, accessed: 2025-04-11 at 1249PM.

12. Ernst, M.D.: Natural Language is a Programming Language: Applying Natural Language Processing to Software Development. In: Lerner, B.S., Bodík, R., Krishnamurthi, S. (eds.) 2nd Summit on Advances in Programming Languages (SNAPL 2017). Leibniz International Proceedings in Informatics (LIPIcs), vol. 71, pp. 4:1–4:14. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl, Germany (2017). https://doi.org/10.4230/LIPIcs.SNAPL.2017.4, https://drops.dagstuhl.de/entities/document/10.4230/LIPIcs.SNAPL.2017.4

13. Fan, W., Rego, M., Hu, X., Dod, S., Ni, Z., Xie, D., DiVincenzo, J., Tan, L.: Evaluating the ability of large language models to generate verifiable specifications in verifast (2025), https://arxiv.org/abs/2411.02318

14. Fang, W., Li, M., Li, M., Yan, Z., Liu, S., Zhang, H., Xie, Z.: Assertllm: Generating hardware verification assertions from design specifications via multi-llms. In: 2024 IEEE LLM Aided Design Workshop (LAD). pp. 1–1 (2024). https://doi.org/10.1109/LAD62341.2024.10691792

15. Fazelnia, M., Mirakhorli, M., Bagheri, H.: Translation titans, reasoning challenges: Satisfiability-aided language models for detecting conflicting requirements. In: Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering. p. 2294–2298. ASE '24, Association for Computing Machinery, New York, NY, USA (2024). https://doi.org/10.1145/3691620.3695302, https://doi.org/10.1145/3691620.3695302

16. Gervasi, V., Nuseibeh, B.: Lightweight validation of natural language requirements. Softw. Pract. Exper. **32**(2), 113–133 (Feb 2002). https://doi.org/10.1002/spe.430, https://doi.org/10.1002/spe.430

17. Ghosh, S., Elenius, D., Li, W., Lincoln, P., Shankar, N., Steiner, W.: Arsenal: Automatic requirements specification extraction from natural language. In: Rayadurgam, S., Tkachuk, O. (eds.) NASA Formal Methods. pp. 41–46. Springer International Publishing, Cham (2016)

18. Granberry, G., Ahrendt, W., Johansson, M.: Specify what? enhancing neural specification synthesis by symbolic methods. In: Kosmatov, N., Kovács, L. (eds.) Integrated Formal Methods. pp. 307–325. Springer Nature Switzerland, Cham (2025)

19. Granberry, G., Ahrendt, W., Johansson, M.: Towards integrating copiloting and formal methods. In: Margaria, T., Steffen, B. (eds.) Leveraging Applications of Formal Methods, Verification and Validation. Specification and Verification. pp. 144–158. Springer Nature Switzerland, Cham (2025)

20. Greenspan, S.J., Borgida, A., Mylopoulos, J.: A requirements modeling language and its logic. Information Systems **11**(1), 9–23 (1986). https://doi.org/https://doi.org/10.1016/0306-4379(86)90020-7, https://www.sciencedirect.com/science/article/pii/0306437986900207

21. Hahn, C., Schmitt, F., Tillman, J.J., Metzger, N., Siber, J., Finkbeiner, B.: Formal specifications from natural language (2022), https://arxiv.org/abs/2206.01962

22. Hsieh, C., Chuang, Y., Li, C., Wang, Z., Le, L.T., Kumar, A., Glass, J.R., Ratner, A., Lee, C., Krishna, R., Pfister, T.: Found in the middle: Calibrating positional attention bias improves long context utilization. In: Ku, L., Martins, A., Srikumar, V. (eds.) Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024. pp. 14982–14995. Association for Computational Linguistics (2024). https://doi.org/10.18653/V1/2024.FINDINGS-ACL.890, https://doi.org/10.18653/v1/2024.findings-acl.890

23. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. In: The Tenth International Conference on Learning Represen-

tations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net (2022), https://openreview.net/forum?id=nZeVKeeFYf9

24. Huisman, M., Gurov, D., Malkis, A.: Formal methods: From academia to industrial practice. a travel guide (2024), https://arxiv.org/abs/2002.07279

25. Jiang, A.Q., Li, W., Tworkowski, S., Czechowski, K., Odrzygóźdź, T., Mił oś, P., Wu, Y., Jamnik, M.: Thor: Wielding hammers to integrate language models and automated theorem provers. In: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (eds.) Advances in Neural Information Processing Systems. vol. 35, pp. 8360–8373. Curran Associates, Inc. (2022)

26. Kojima, T., Gu, S.S., Reid, M., Matsuo, Y., Iwasawa, Y.: Large language models are zero-shot reasoners. In: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (eds.) Advances in Neural Information Processing Systems. vol. 35, pp. 22199–22213. Curran Associates, Inc. (2022)

27. Leong, I.T., Barbosa, R.: Translating natural language requirements to formal specifications: A study on gpt and symbolic nlp. In: 2023 53rd Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W). pp. 259–262 (2023). https://doi.org/10.1109/DSN-W58399.2023.00065

28. Leopold, H., Mendling, J., Polyvyanyy, A.: Supporting process model validation through natural language generation. IEEE Transactions on Software Engineering **40**(8), 818–840 (2014). https://doi.org/10.1109/TSE.2014.2327044

29. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W., Rocktäschel, T., Riedel, S., Kiela, D.: Retrieval-augmented generation for knowledge-intensive NLP tasks. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual (2020), https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html

30. Li, J., Li, G., Li, Y., Jin, Z.: Structured chain-of-thought prompting for code generation. ACM Trans. Softw. Eng. Methodol. **34**(2) (Jan 2025). https://doi.org/10.1145/3690635, https://doi.org/10.1145/3690635

31. Li, M., Fang, W., Zhang, Q., Xie, Z.: Specllm: Exploring generation and review of vlsi design specification with large language model (2024), https://arxiv.org/abs/2401.13266

32. Li, Y., Hui, B., Xia, X., Yang, J., Yang, M., Zhang, L., Si, S., Chen, L.H., Liu, J., Liu, T., Huang, F., Li, Y.: One-shot learning as instruction data prospector for large language models (2024), https://arxiv.org/abs/2312.10302

33. Luckcuck, M., Farrell, M., Dennis, L.A., Dixon, C., Fisher, M.: Formal specification and verification of autonomous robotic systems: A survey. ACM Comput. Surv. **52**(5) (Sep 2019). https://doi.org/10.1145/3342355, https://doi.org/10.1145/3342355

34. Ma, L., Liu, S., Li, Y., Xie, X., Bu, L.: Specgen: Automated generation of formal program specifications via large language models (2024), https://arxiv.org/abs/2401.08807

35. Mandal, S., Chethan, A., Janfaza, V., Mahmud, S.M.F., Anderson, T.A., Turek, J., Tithi, J.J., Muzahid, A.: Large language models based automatic synthesis of software specifications (2023), https://arxiv.org/abs/2304.09181

36. Misu, M.R.H., Lopes, C.V., Ma, I., Noble, J.: Towards ai-assisted synthesis of verified dafny methods. Proc. ACM Softw. Eng. **1**(FSE) (Jul 2024). https://doi.org/10.1145/3643763, https://doi.org/10.1145/3643763

37. Mugnier, E., Gonzalez, E.A., Jhala, R., Polikarpova, N., Zhou, Y.: Laurel: Generating dafny assertions using large language models (2024), https://arxiv.org/abs/2405.16792

38. Mukherjee, P., Delaware, B.: Towards automated verification of llm-synthesized c programs (2024), https://arxiv.org/abs/2410.14835

39. Nayak, A., Timmapathini, H.P., Murali, V., Ponnalagu, K., Venkoparao, V.G., Post, A.: Req2spec: Transforming software requirements into formal specifications using natural language processing. In: Requirements Engineering: Foundation for Software Quality: 28th International Working Conference,

REFSQ 2022, Birmingham, UK, March 21–24, 2022, Proceedings. p. 87–95. Springer-Verlag, Berlin, Heidelberg (2022)

40. Nazaruka, E., Osis, J.: Determination of natural language processing tasks and tools for topological functioning modelling. In: Proceedings of the 13th International Conference on Evaluation of Novel Approaches to Software Engineering. p. 501–512. ENASE 2018, SCITEPRESS - Science and Technology Publications, Lda, Setubal, PRT (2018). https://doi.org/10.5220/0006817205010512, https://doi.org/10.5220/0006817205010512

41. Necula, S.C., Dumitriu, F., Greavu-Șerban, V.: A systematic literature review on using natural language processing in software requirements engineering. Electronics **13**(11) (2024). https://doi.org/10.3390/electronics13112055, https://www.mdpi.com/2079-9292/13/11/2055

42. Nelken, R., Francez, N.: Automatic translation of natural language system specifications into temporal logic. In: Alur, R., Henzinger, T.A. (eds.) Computer Aided Verification. pp. 360–371. Springer Berlin Heidelberg, Berlin, Heidelberg (1996)

43. Nouri, A., Cabrero-Daniel, B., Törner, F., Sivencrona, H., Berger, C.: Engineering safety requirements for autonomous driving with large language models. In: 2024 IEEE 32nd International Requirements Engineering Conference (RE). pp. 218–228 (2024). https://doi.org/10.1109/RE59067.2024.00029

44. Nowakowski, W., Śmiałek, M., Ambroziewicz, A., Straszak, T.: Requirements-level language and tools for capturing software system essence. Computer Science and Information Systems **10**(4), 1499–1524 (2013)

45. Osborne, M., MacNish, C.: Processing natural language software requirement specifications. In: Proceedings of the Second International Conference on Requirements Engineering. pp. 229–236 (1996). https://doi.org/10.1109/ICRE.1996.491451

46. Palin, R., Ward, D., Habli, I., Rivett, R.: Iso 26262 safety cases: Compliance and assurance. In: 6th IET International Conference on System Safety 2011. pp. 1–6 (2011). https://doi.org/10.1049/cp.2011.0251

47. Preda, A.R., Mayr-Dorn, C., Mashkoor, A., Egyed, A.: Supporting high-level to low-level requirements coverage reviewing with large language models. In: Proceedings of the 21st International Conference on Mining Software Repositories. p. 242–253. MSR '24, Association for Computing Machinery, New York, NY, USA (2024). https://doi.org/10.1145/3643991.3644922, https://doi.org/10.1145/3643991.3644922

48. Quan, X., Valentino, M., Dennis, L.A., Freitas, A.: Verification and refinement of natural language explanations through llm-symbolic theorem proving (2024), https://arxiv.org/abs/2405.01379

49. Reinpold, L.M., Schieseck, M., Wagner, L.P., Gehlhoff, F., Fay, A.: Exploring llms for verifying technical system specifications against requirements (2024), https://arxiv.org/abs/2411.11582

50. Shum, K., Diao, S., Zhang, T.: Automatic prompt augmentation and selection with chain-of-thought from labeled data. In: Bouamor, H., Pino, J., Bali, K. (eds.) Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023. pp. 12113–12139. Association for Computational Linguistics (2023). https://doi.org/10.18653/V1/2023.FINDINGS-EMNLP.811, https://doi.org/10.18653/v1/2023.findings-emnlp.811

51. Tihanyi, N., Jain, R., Charalambous, Y., Ferrag, M.A., Sun, Y., Cordeiro, L.C.: A new era in software security: Towards self-healing software via large language models and formal verification (2024), https://arxiv.org/abs/2305.14752

52. Wei, J., Wang, X., Schuurmans, D., Bosma, M., ichter, b., Xia, F., Chi, E., Le, Q.V., Zhou, D.: Chain-of-thought prompting elicits reasoning in large language models. In: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (eds.) Advances in Neural Information Processing Systems. vol. 35, pp. 24824–24837. Curran Associates, Inc. (2022)

53. White, J., Fu, Q., Hays, S., Sandborn, M., Olea, C., Gilbert, H., Elnashar, A., Spencer-Smith, J., Schmidt, D.C.: A prompt pattern catalog to enhance prompt engineering with chatgpt. CoRR **abs/2302.11382** (2023). https://doi.org/10.48550/ARXIV.2302.11382, https://doi.org/10.48550/arXiv.2302.11382

54. Wu, H., Barrett, C., Narodytska, N.: Lemur: Integrating large language models in automated program verification (2024), https://arxiv.org/abs/2310.04870
55. Xu, W., Banburski-Fahey, A., Jojic, N.: Reprompting: Automated chain-of-thought prompt inference through gibbs sampling. CoRR **abs/2305.09993** (2023). https://doi.org/10.48550/ARXIV.2305.09993, https://doi.org/10.48550/arXiv.2305.09993
56. Xu, Y., Feng, J., Miao, W.: Learning from failures: Translation of natural language requirements into linear temporal logic with large language models. In: 2024 IEEE 24th International Conference on Software Quality, Reliability and Security (QRS). pp. 204–215 (2024). https://doi.org/10.1109/QRS62785.2024.00029
57. Yan, R., Cheng, C.H., Chai, Y.: Formal consistency checking over specifications in natural languages. In: 2015 Design, Automation & Test in Europe Conference & Exhibition (DATE). pp. 1677–1682 (2015)
58. Yang, K., Swope, A., Gu, A., Chalamala, R., Song, P., Yu, S., Godil, S., Prenger, R.J., Anandkumar, A.: Leandojo: Theorem proving with retrieval-augmented language models. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (eds.) Advances in Neural Information Processing Systems. vol. 36, pp. 21573–21612. Curran Associates, Inc. (2023)
59. Yao, J., Liu, Y., Dong, Z., Guo, M., Hu, H., Keutzer, K., Du, L., Zhou, D., Zhang, S.: Promptcot: Align prompt distribution via adapted chain-of-thought. In: 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 7027–7037 (2024). https://doi.org/10.1109/CVPR52733.2024.00671
60. Ye, X., Durrett, G.: The unreliability of explanations in few-shot prompting for textual reasoning. In: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (eds.) Advances in Neural Information Processing Systems. vol. 35, pp. 30378–30392. Curran Associates, Inc. (2022)
61. Zhang, H., Cai, M., Zhang, X., Zhang, C.J., Mao, R., Wu, K.: Self-convinced prompting: Few-shot question answering with repeated introspection (2023), https://arxiv.org/abs/2310.05035

# 8   Appendices

## 8.1   Dafny Lemma Definition

The Dafny lemma, `ParseDigitsAndDot`, specifies the expected behaviour of the function `ParseDecStr`. It operates on an input string composed of digits and a decimal separator.

**Lemma Definition:**

```
lemma ParseDigitsAndDot(s1: string, s2: string)
  requires ∀ i | 0 ≤ i < |s1| :: '0' ≤ s1[i] ≤ '9'
  ensures ParseDecStr(s1+"."+s2).value.1 == "."+s2
  {
    if |s1| == 1 {
      assert ParseDecStr("."+s2).None?;
    } else {
      ParseDigitsAndDot(s1[1..],s2);
      assert s1 + "." + s2 == [s1[0]] + (s1[1..] + "." + s2);
    }
  }
```

## 8.2   Summary Tables

| Ref. | Tool / Framework / Technique | Description |
|---|---|---|
| [10] | LLM-based Code Verification | Uses LLMs like GPT-3.5 to verify code by analyzing requirements and explaining whether they are met. |
| [9] | nl2spec | A framework leveraging LLMs to generate formal specifications from natural language, addressing ambiguity in system requirements with iterative refinement. |
| [48] | Explanation-Refiner | A neuro-symbolic framework integrating LLMs and theorem provers to formalize and validate explanatory sentences, providing error correction and feedback for improving NLI models. |
| [27] | *Not Specified* | Analyzes research directions in software requirement engineering, conducting a SWOT analysis and sharing evaluation findings. |
| [3] | Symbolic NLP vs. ChatGPT | Compares the performance of symbolic NLP and ChatGPT in generating correct JML output from natural language preconditions. |
| [2] | Domain Model Extractor | Generates domain models from natural language requirements in an industrial case study, evaluating accuracy and performance. |
| [35] | SpecSyn | A framework using LLMs for automatic synthesis of software specifications, improving accuracy by 21% over previous tools. |
| [14] | AssertLLM | A tool generating assertions for hardware verification from design specifications using three customized LLMs, achieving 89% correctness. |
| [16] | Formal Verification of NASA's Software | Reports on formal verification of NASA's Node Control Software natural language specifications, highlighting errors and lessons learned. |
| [31] | SpecLLM | Explores using LLMs for generating and reviewing VLSI design specifications, improving chip design documentation. |
| [44] | Requirements Specification Language (RSL), ReDSeeDS | Enhanced software requirements specification using constrained natural language and automated transformations into code. |
| [17] | ARSENAL Framework and Methodology | Automated extraction of requirements specification from natural language with automatic verification. |
| [28] | BPM-to-NL Translation Process | Generated natural language descriptions from business process models for better validation. |
| [37] | Laurel | A framework to generate Dafny assertions to automate program verification process for a SMT solver |
| [45] | Controlled Natural Language (CL) with ANLT | Expressed software requirements in a limited set of natural language and translated to logical expressions to detect ambiguities. |
| [49] | LLM-based Analysis for Smart Grid Requirements | Improved smart grid requirement specifications with GPT-4o and Claude 3.5 Sonnet, achieving F1-scores between 79% - 94%. |
| [56] | NL-to-LTL Translation via LLMs | Converted unstructured natural language requirements to NL-LTL pairs, achieving 94.4% accuracy on public datasets. |
| [51] | ESBMC-AI | Combined LLMs with Formal Verification to detect and fix software vulnerabilities with high accuracy. |
| [21] | LLM-based Formal Specifications Translation | Translated natural language into formal rules (regex, FOL, LTL) with high adaptability and performance. |
| [38] | SynVer Framework | Synthesized and verified C programs using the Verified Software Toolchain. |
| [47] | LLM-based Requirement Coverage Analysis | Ensured low-level software requirements met high-level requirements, achieving 99.7% recall in spotting missing coverage. |
| [15] | SAT-LLM | Integrated SMT solvers with LLMs to improve conflict identification in requirements; significantly outperformed standalone LLMs in detecting complex conflicts. |

Table 2: Summary of LLMs related literature (Part 1 of 2)

| Ref. | Tool / Framework / Technique | Description |
|---|---|---|
| [12] | NLP for Software Development | Assessed NLP techniques for various software development stages, highlighting their suitability for generating assertions and processing developer queries. |
| [39] | Req2Spec | NLP-based tool that formalises natural language requirements for HANFOR; achieved 71% accuracy in formalising 222 automotive requirements at BOSCH. |
| [20] | RML (Requirements Modeling Language) | Introduced a conceptual model-based framework ensuring precision, consistency, and clarity in requirements writing. |
| [13] | GPT-4o for VeriFast Verification | Evaluated GPT-4o's ability to generate C program specifications for VeriFast; found that while functional behavior was preserved, verification often failed or contained redundancies. |
| [42] | NL to Temporal Logic Translation | Developed an automatic translation mechanism from natural language sentences to temporal logic for formal verification. |
| [8] | ANTONIO toolkit | A comprehensive analysis of NLP verification approaches and introduces a structured **NLP Verification Pipeline** with six key components. The work includes identifying gaps in existing methods, proposing novel solutions for improved robustness, extending standard verifiability metrics, and emphasizing the importance of reporting verified subspace (geometric and semantic) properties for better reliability and interpretability. |
| [54] | Lemur | Integrated LLMs with automated reasoners for program verification, defining sound transition rules and demonstrating improved performance on benchmark tests. |
| [41] | Systematic Review | Conducted a comprehensive review on natural language to formal specification translation, analyzing research across multiple academic databases. |
| [43] | LLM-based Safety Requirements Pipeline | Designed a pipeline using LLMs to refine and decompose safety requirements for autonomous vehicles, evaluated through expert assessments and industrial implementation. |
| [57] | Specification Consistency Framework | Ensured consistency between oral and formal specifications, incorporating time extraction, input-output partitioning, and semantic reasoning, with positive evaluation results. |
| [40] | NLP Tools for TFM | Evaluated six NLP pipelines for Topological Functioning Modelling (TFM). Found that Stanford CoreNLP, FreeLing, and NLTK performed best. |
| [36] | LLM-based Dafny Task Generation | Used LLMs (GPT-4, PaLM-2) to generate Dafny tasks from MBPP benchmark using different prompting strategies (context-less, signature, retrieval-augmented CoT). GPT-4 achieved best results with retrieval-augmented CoT prompt, producing 153 verified Dafny solutions. |
| [58] | LeanDojo & ReProver | Introduced LeanDojo, an open-source toolkit for interacting with the Lean theorem prover. Developed ReProver, a retrieval-augmented LLM-based prover that improved theorem proving efficiency. Created a benchmark with 98,734 theorems and proofs for testing generalization. |
| [25] | Thor and class methods named Hammers | Introduced a framework named Thor which integrates language models with theorem provers. Hammers are implemented to find the appropriate premises to complete the proofs of conjectures. Datasets used are PISA and MiniF2F. |
| [19] | Not specified | The work proposes the integration of major formal languages (Dafny, Ada/SPARK, Frama-C, and KeY), their interactive theorem provers (Coq, Isabelle/HOL, Lean) with Copilot. |
| [33] | Systematic Review | Conducted a comprehensive survey on formal specification and verification of autonomous robotic systems in 2018. This is based on literature available of ten years (2008 - 2018). |
| [18] | Symbolic analysis and LLMs prompts | The quality of annotations produced in ACSL format is measured for PathCrawler and EVA (tools available in Frama-C). PathCrawler generated more context-aware annotations while, EVA efficiency improved having less run-time errors. |

Table 3: Summary of LLMs related literature (Part 2 of 2)

Table 4: Classification and Description of Surveyed Literature

| Ref. | Tool / Work | Classification | Description |
|---|---|---|---|
| [2] | Domain Model Extractor | Fine-tuned | Generates domain models from NL requirements; evaluated in an industrial case study for performance and accuracy. |
| [3] | Symbolic NLP vs. ChatGPT | Prompt-only | Compares symbolic NLP and ChatGPT in generating correct JML from NL preconditions. |
| [28] | BPM-to-NL Translation | Prompt-only | Translates business process models to NL to support better stakeholder validation. |
| [54] | Lemur | Verifier-in-loop | Integrates LLMs with automated reasoners and defines sound transition rules for verification. |
| [13] | GPT-4o for VeriFast Verification | Verifier-in-loop | Assesses GPT-4o's performance in generating C specs for VeriFast; captures issues in functional correctness. |
| [17] | ARSENAL | Fine-tuned | Extracts requirements from NL and performs automatic verification. |
| [37] | Laurel for Dafny | Prompt-only + Verifier-in-loop | Automates generation of Dafny assertions to support SMT-based verification. |
| [18] | PathCrawler + EVA | Verifier-in-loop + Prompt | PathCrawler generates context-aware ACSL annotations; EVA reduces runtime errors in Frama-C. |
| [42] | NL to Temporal Logic Translation | Prompt-only | Automatically translates NL into temporal logic for formal verification. |
| [39] | Req2Spec | Prompt-only | Converts NL requirements into formal specs (e.g., HANFOR); 71% accuracy on BOSCH data. |
| [14] | AssertLLM | Multi-LLMs / Prompt-only | Uses 3 customized LLMs to generate assertions from hardware design specs; 89% correctness achieved. |
| [35] | SpecSyn | Fine-tuned | Synthesizes software specifications from NL, improving over prior tools by 21%. |
| [9] | nl2spec | Prompt-only + Iterative Refinement | Iteratively generates formal specs from NL requirements, reducing ambiguity. |
| [47] | LLM-based Requirement Coverage | Prompt-only | Maps low-level requirements to high-level ones with 99.7% recall in coverage detection. |
| [31] | SpecLLM | Prompt-only | Uses LLMs to create and review VLSI design specs, enhancing chip documentation. |
| [56] | NL-to-LTL via LLMs | Prompt-only | Converts unstructured NL to NL-LTL pairs with 94.4% accuracy. |
| [8] | ANTONIO Toolkit | Verifier-in-loop | Introduces an NLP Verification Pipeline with metrics, gaps, and semantic subspace verification proposals. |
| [10] | GPT-3.5 for Code Verification | Prompt-only | Uses GPT-3.5 to verify code against requirements, providing feedback on requirement satisfaction. |
| [49] | GPT-4o + Claude for Smart Grid | Verifier-in-loop | Applies GPT-4o and Claude 3.5 for smart grid requirement verification, reaching 79–94% F1-scores. |
| [41] | Systematic Review | Meta-analysis | Surveys NL-to-specification literature across domains and academic sources. |
| [15] | SAT-LLM | Neuro-symbolic (LLM + SMT) | Combines LLMs with SMT to detect complex conflicts in requirements. |
| [25] | Thor | Neuro-symbolic | Integrates LLMs with theorem provers using class methods like Hammers for proof completion. |
| [36] | Dafny Task Gen w/ CoT | Prompt + Retrieval + CoT | GPT-4 and PaLM-2 generate verified Dafny tasks via retrieval-augmented CoT prompting. |
| [58] | LeanDojo + ReProver | Retrieval-augmented | Retrieval-augmented LLM-based prover improves Lean theorem proving on 98K+ samples. |
| [38] | SynVer for C | Verifier-in-loop | Synthesizes and verifies C programs using VST with improved automation. |
| [19] | Copilot + Formal Methods | IDE Integration Proposal | Suggests integrating formal tools (e.g., Dafny, Coq, KeY) into IDEs like Copilot. |
| [33] | Robotic Systems Review | Meta-analysis | Reviews 10 years of literature on formal verification in autonomous robotic systems. |
| [12] | NLP for Software Dev | Survey / Meta-analysis | Evaluates NLP techniques in software development life cycle. |
| [45] | RML (1986) | Manual / Controlled NL | Introduces controlled NL framework for precise and consistent requirements writing. |
| [51] | ESBMC-AI | Neuro-symbolic | Uses LLMs + formal verification to detect vulnerabilities in software. |
| [57] | Specification Consistency Framework | Manual / Baseline | Aligns oral and formal specifications using semantic reasoning and input-output analysis. |
| [43] | LLM Safety Req. Pipeline | Prompt-only | Uses LLMs to decompose and refine autonomous vehicle safety requirements. |
| [40] | NLP Tools for TFM | Prompt-only | Compares NLP pipelines for topological modeling; CoreNLP and FreeLing perform best. |

Table 5: Classification of Surveyed Literature by Methodology

| Ref. | Tool / Work | Classification |
|---|---|---|
| [6] | Prompt Engineering Survey | Prompt-only |
| [46] | Advancing RE with LLMs | Prompt-only |
| [4] | Domain Model Extractor | Fine-tuned |
| [24] | IEC 61508 Safety Standard | Not LLM-based |
| [37] | Reasoning Topologies | Prompt-only |
| [11] | DO-178C Standard | Not LLM-based |
| [10] | BIS NL2SQL | Prompt-only |
| [9] | ANTONIO Toolkit | Verifier-in-loop |
| [48] | nl2spec | Prompt-only + Iterative Refinement |
| [2] | GPT-3.5 for Code Verification | Prompt-only |
| [27] | Elicit Tool | Tool Support (Meta) |
| [3] | NLP in Software Dev | Survey / Meta-analysis |
| [35] | GPT-4o + VeriFast | Verifier-in-loop |
| [14] | AssertLLM | Multi-LLMs / Prompt-only |
| [16] | SAT-LLM | Neuro-symbolic (LLM + SMT) |
| [31] | NASA SW Formal Verification | Manual / Baseline |
| [44] | ARSENAL | Fine-tuned |
| [17] | PathCrawler + EVA | Verifier-in-loop + Prompt |
| [28] | Copilot + Formal Methods | IDE Integration Proposal |
| [45] | RML (1986) | Manual / Controlled NL |
| [49] | NL to FOL, Regex, LTL | Prompt-only |
| [56] | Long Context Calibration | Prompt-only |
| [15] | LoRA Tuning | Fine-tuned |
| [12] | Formal Methods Transfer | Meta-analysis |
| [39] | Thor (LLM + Theorem Prover) | Neuro-symbolic |
| [20] | Zero-shot Reasoning | Prompt-only |
| [13] | Symbolic NLP vs ChatGPT | Prompt-only |
| [42] | BPM to NL Translation | Prompt-only |
| [58] | RAG for NLP Tasks | Retrieval-augmented |
| [25] | Structured CoT for Code | Prompt-only |
| [19] | SpecLLM | Prompt-only |
| [18] | One-shot Learning | Prompt-only |
| [34] | Robotic Systems Review | Meta-analysis |
| [7] | SpecGen | Prompt + Mutation (Verifier-in-loop) |
| [26] | SpecSyn | Fine-tuned |
| [32] | Dafny Task Gen w/ CoT | Prompt + Retrieval + CoT |
| [60] | Laurel for Dafny | Prompt-only + Verifier-in-loop |
| [61] | SynVer for C | Verifier-in-loop |
| [22] | Req2Spec | Prompt-only |
| [52] | NLP Tools for TFM | Prompt-only |
| [42] | NL to Temporal Logic | Fine-tuned |
| [22] | LLMs in Autonomous Driving | Prompt-only |
| [44] | RSL + ReDSeeDS | Fine-tuned / Controlled NL |
| [45] | Controlled NL (1996) | Controlled NL |
| [46] | ISO 26262 | Not LLM-based |
| [10] | LLMs for Requirement Coverage | Prompt-only |
| [48] | Explanation-Refiner | Neuro-symbolic |
| [49] | GPT-4o / Claude 3.5 for Smart Grid | Prompt-only |
| [50] | Automate-CoT | Prompt-only + Automated CoT |
| [51] | ESBMC-AI | Verifier-in-loop |
| [52] | Chain-of-Thought (CoT) | Prompt-only |
| [53] | Prompt Pattern Catalog | Prompt-only |
| [54] | Lemur (LLM + Verifier) | Verifier-in-loop |
| [55] | Reprompting w/ Gibbs | Prompt-only |
| [56] | NL to LTL (Dynamic Prompt) | Prompt + Human-in-loop |
| [57] | NL Consistency Framework | Fine-tuned |
| [58] | LeanDojo + ReProver | Retrieval + Verifier-in-loop |
| [59] | PromptCoT | Prompt-only |
| [60] | Explanation Failures in Few-Shot | Prompt-only |
| [61] | Self-Convinced Prompting | Prompt-only |