

Bridging Natural Language and Code by Transforming Free-Form Sentences into Sequence of Unambiguous Sentences with Large Language Model

Nikita Kiran Yeole

Computer Science

Virginia Tech, Blacksburg, USA

nikitay@vt.edu

Michael S. Hsiao

Electrical and Computer Engineering

Virginia Tech, Blacksburg, USA

hsiao@vt.edu

Abstract—In the realm of natural language programming, translating free-form sentences in natural language into a functional, machine-executable program remains difficult due to the following 4 challenges. First, the inherent ambiguity of natural languages. Second, the high-level verbose nature in user descriptions. Third, the complexity in the sentences and Fourth, the invalid or semantically unclear sentences. Our proposed solution is a Large language model based Artificial Intelligence driven assistant to process free-form sentences and decompose them into sequences of simplified, unambiguous sentences that abide by a set of rules, thereby stripping away the complexities embedded within the original sentences. These resulting sentences are then used to generate the code. We applied the proposed approach to a set of free-form sentences written by middle-school students for describing the logic behind video games. More than 60 percent of the free-form sentences containing these problems were successfully converted to sequences of simple unambiguous object-oriented sentences by our approach.

Keywords—Natural language programming; decomposition; chain-of-thought reasoning.

I. INTRODUCTION

Natural Language Programming (NLPg) is a concept that attempts to convert instructions/specifications written in free-form natural language into functional program code. NLPg envisions a world in which everyone can program machines without understanding the intricacies of conventional programming languages. While generative Artificial Intelligence (AI) has shown some success in producing code snippets from natural language text, the code that is produced may not adhere to the intent of the input text. When the code does not meet the intent, the user can do one of two things: (1) manually modify the generated code, or (2) re-write the natural language text and try to generate new code. For users who are not experienced programmers, option 1 may not be feasible, since the generated code may contain data structures and/or algorithms that the user is unfamiliar with. Hence, the user is left with the second option. In order to generate functionally correct code, the input must be in a format that the system can process such that common problems with general natural languages are removed. In other words, if the input text is semantically unambiguous, the code generated will more likely adhere to the intent of the input text [1].

An additional benefit is that this helps the user to learn to write unambiguous input text, a necessary skill behind the thought processes in coding. Natural language is increasingly

applied in education for personalized AI tutoring and interactive learning, aiding educators in various ways [2] [3] [4]. The ability to instruct a machine in natural language bridges the gap between human thought processes and the digital world, making technology more accessible and intuitive for students.

There are many factors associated with natural language instructions, which makes NLPg extremely challenging [5]. First, the ambiguity in the sentences. Second, the high level verbose descriptions given by humans. Third, complex and compound sentences. Fourth, invalid or erroneous sentences written by humans. We will briefly highlight each of these four areas in the following discussion.

Natural Language (NL) sentences can include ambiguities wherein a single word or phrase may have several interpretations. Consider, for instance, the following English sentence employed in game design:

"When the rabbit touches a rock, it explodes."

Here, the phrase containing the pronoun 'it' creates uncertainty in this sentence. According to one view, the rabbit explodes after touching the rock, whereas the other contends that the rock explodes.

Secondly, the NL instructions can be excessively verbose, especially written by the people who may not know how to program. Consider, for instance, the English sentence employed in game design:

"In a mysterious realm, a lone pointer and some aliens engage in a cosmic dance. When the pointer touches an alien, it changes colors: original to purple, purple to pink. Pink aliens explode."

Here, the sentences provided are verbose with extraneous descriptive words and phrases. Although they adhere to proper English grammar, they deviate from a concise format.

Thirdly, machines typically demand sentences with a clear structure containing a subject, verb, and object. However, complex sentences that sequentially combine multiple events may complicate the parsing of the sentence and prevent a full understanding of the intent of the user. The following sentence illustrates one such example:

"When the carrot turns into a diamond before the carrot touches a fox, the score increases."

Fourthly, when humans provide instructions, there is a chance that they might offer sentences that are invalid, illogical, incomplete or erroneous. In such cases, it becomes

difficult for the machine to extract the exact task that needs to be executed. The following is one such example:

"Brick spawns at the bottom. 14 cheese at the top in rows. Ball in the middle. w is up. s is down. brick touches border bounce. ball touches cheese bounces back."

To overcome these challenges, we propose an Artificial Intelligence driven assistant using Large Language Models (LLMs), which will attempt to convert the free-form sentences into sequences of simple sentences, each with a clear subject, verb, and object structure. It promotes a paradigm where instead of the user conforming to the machine, the machine adapts to grasp the user's intent. This assistant streamlines, simplifies, and transforms the NL phrases into directives that machines can easily interpret. The design of the assistant prioritizes rule-driven simplification, methodically translating sentences that eliminate unnecessary elements while retaining the core meaning.

Motivating Example: Consider the following free-form description of a game:

"The rabbit wanders, reversing at borders. The fox wanders, chasing the rabbit when spotting the rabbit. When the rabbit touches the fox, the fox turns into a carrot."

Our goal is to convert the above paragraph to the following simplified sentences.

"There is a rabbit. There is a fox. The rabbit wanders. The fox wanders. If the rabbit reaches a border, it reverses. If the fox sees the rabbit, it chases the rabbit. When the rabbit touches the fox, the fox becomes mutated. When the fox is mutated, it turns into a carrot."

The deconstruction of complex sentences and then rewriting them in basic, simple sentences is the most novel aspect of our strategy. The NL expression frequently combines various thoughts or directives in a single, complex sentence [6]. So, these sentences are decomposed and rewritten in a format that abides by imposed rules. In our approach, the input sentences are parsed, during which the engine identifies key components and breaks them down into their basic elements. By analyzing the relationships between these elements, the system deciphers the user's intention. With this insight, it reconstructs the information into simple sentences that are structured and guided by rules.

The novelty of this paper lies in its specific methodology for simplifying natural language sentences into structured directives through a rule-based system, a departure from traditional semantic parsing and tree-based neural network models which often struggle with the ambiguity and complexity of natural language [5]. We also integrate an educational platform, GameChangineer, to demonstrate the practical application of this approach, showcasing how it facilitates the learning of object-oriented programming concepts by converting these simplified sentences into functional game code.

We applied our approach to process 1000 free-write sentences, out of which 800 sentences contained at least one of the four aforementioned problems, and 200 sentences are non-problematic sentences. The rewritten sentences are then given to an educational platform called GameChangineer [7], [8]

that can convert the object oriented English sentences to a functional game [9]. GameChangineer is an AI-Enabled Design and Education Platform which helps students to discover and practice logical reasoning, problem-solving, algorithmic design, critical and computational thinking [7]. Beginners may find Object-Oriented Programming (OOP) to be abstract and challenging to understand due to its emphasis on classes, objects, inheritance, polymorphism, encapsulation, and abstraction. Students can express their thoughts and queries in a way that comes naturally to them when they are able to interact with an educational software through natural language. This reduces the cognitive load associated with learning new, technical syntax and concepts, allowing them to focus more on the underlying principles of OOP. The results showed that more than 60% of the problematic sentences were successfully converted by our approach. The sentences which were successfully converted led to a correct, functional game which adheres to the intent of the user.

The rest of the paper is organized as follows. Section 2 describes the related work. Section 3 lays out the methodology in our work and Section 4 presents the evaluation of our approach and discusses its implications. Finally, Section 5 concludes the paper.

II. RELATED WORK

A curated list of groundbreaking studies that have had an impact on this field is included in this section.

One approach to addressing these natural language challenges is through semantic parsing, where natural language utterances are encoded and translated into syntactically correct target code snippets using tree-based neural network models [5]. This technique shows promise in generating accurate code snippets from natural language descriptions by focusing on the structural aspects of language to reduce ambiguity and manage complexity. Even sophisticated semantic parsing models, while capable of generating syntactically correct code from natural language inputs, often face difficulties in capturing the user's intent accurately. This is because a single phrase can be interpreted in multiple ways, leading to code that, while technically correct, does not fulfill the intended function [5].

Another sophisticated method involves using execution-based selection processes and Minimum Bayes Risk (MBR) decoding to minimize expected errors in the generated code [10]. This approach selects the most accurate output by considering the execution results of the generated code samples, helping to ensure that the generated code aligns with the intended functionality described in natural language. This approach has its limitations. It requires executing several generated code snippets to determine the best candidate, which can be computationally expensive and inefficient. Furthermore, if the initial pool of generated code contains errors or fails to capture the user's intent accurately, the selection process may still result in sub-optimal code [10].

Deep learning techniques offer significant advancements in understanding and generating code from natural language. By leveraging the encoder-decoder framework, these models can

learn from vast datasets of code to improve the accuracy and relevance of generated code snippets, addressing issues of verbosity and complex sentence structures by focusing on the semantic content of the instructions [11]. Although deep learning has shown promise in understanding and generating code, the models still struggle with sentences that contain multiple actions or intertwined concepts, reflecting a gap in handling real-world complexity [11]. These limitations underline the necessity for a proposed solution that addresses these core issues.

The Transformer model was first presented by Vaswani et al. in their landmark study, "Attention Is All You Need" [12]. In order to deal with ambiguity, the architecture's self-attention mechanism, which is skilled at capturing context, is essential.

Generative pre-trained transformer (GPT)-3 showed its skill in deciphering a wide range of human expressions and offered a solution to unclear or lacking instructions [13]. Despite its outstanding powers, GPT-3 occasionally produces overly detailed or irrelevant answers [13]. GPT-3 also frequently requires particular fine-tuning for certain tasks [13]. BERT's (Bidirectional Encoder Representations from Transformers) pre-training procedure was improved by Liu et al., who published "RoBERTa: A Robustly Optimized BERT Pretraining Approach" [14] [15].

Wei et al.'s study on "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models" forms a crucial basis for understanding how Chain of Thought (CoT) in LLMs (Large Language Models) can decompose complex reasoning tasks into a series of simpler, logical steps [16]. The authors demonstrate that CoT prompting significantly improves the ability of LLMs to perform complex reasoning tasks across various domains. We employ CoT not for general reasoning enhancement, but specifically for tackling linguistic challenges in programming, such as verbosity, ambiguities, and complex phrase structures.

We focus on preserving the fundamental semantic meaning of the given instructions while simultaneously addressing the inherent difficulties and limitations of human language. The subtleties of freely written phrases can have a profound impact on the semantic meaning, which is the fundamental core of a communication [17]. Therefore, a major goal in this area should be to transform these statements into more straightforward forms without distorting or losing the original meaning that the user intended. This balance makes sure that, despite the language being more structured or standardized for computational processing, the converted sentences remain true to the message the user intended to convey.

III. METHODOLOGY

The foundation of our research is a representative dataset, which was used as the LLM's main input. The data included 1000 student-written free-form sentences as game descriptions. 800 of these sentences have been identified as potentially problematic and 200 sentences have been identified as non-problematic. These descriptions offered a variety of linguistic patterns and semantic complexities. The game descriptions

were diverse, varied in their lengths, and offered a number of difficulties. These sentences showed some ambiguity because they frequently contained intricate structures and relationships that were not always clear. This dataset was also chosen to evaluate the LLM's capacity to comprehend and translate the ambiguous and complex texts into more rule-based, simplified formats.

We used the GPT-3.5 Turbo, a powerful language model created by OpenAI, for the purposes of this study. We made this choice after carefully comparing the performance of GPT-3.5 Turbo and GPT-4, two recent revisions of OpenAI's generative models. Although GPT-4 is a more recent model and is anticipated to offer higher capabilities in many contexts [18], GPT-3.5 Turbo showed improved sentence construction in the most basic form and coherence for the particular prompt utilized in this research. This underscored the need of selecting a model that is tailored to the precise specifications of the work at hand as opposed to just selecting the most recent version. This model was deployed by means of direct integration with the OpenAI API, which allowed us to operate the model locally in our computational environment. Python was selected as our primary programming language because of its extensive libraries for data manipulation and its seamless integration with the OpenAI API.

The model's temperature was set to zero. The choice was made to guarantee deterministic performance from the model.

The top_p parameter was set to 1. This implies that at each stage of the generation process, the model will only take into account the tokens that are the most likely.

It should be emphasized that these combinations signify that we used the model outside of its intended parameters. We purposefully restricted the model to create consistent and repeatable results customized to our needs rather than utilizing its potential for creative and varied outputs. These settings came in helpful in situations where consistency and predictability were crucial.

Our method employed a split strategy that made use of both user prompts and system prompts. The user prompt constitutes the primary interaction point with the user. It is necessary to convert these user-provided free-form sentences into a (sequence of) more simplified structure. The model must understand these inputs robustly due to the inherent variation in how users phrase their queries or utterances. Free-form phrases can be anything from simple sentences to more complex thoughts or assertions, and the challenge lies in distilling the essence of what the user wants to communicate and converting it into a form that the model can process efficiently.

The system prompt serves primarily as a tool to direct the model towards a specific context or mode of operation. We directed the model's potential and ensured that we receive the desired output by creating a structured system prompt. It encompasses a chain-of-thought reasoning via (1) Question Answering, (2) Sentence Reframing, (3) Sentence Decomposition. Figure 1 shows the process flow with an example prompt for each step.

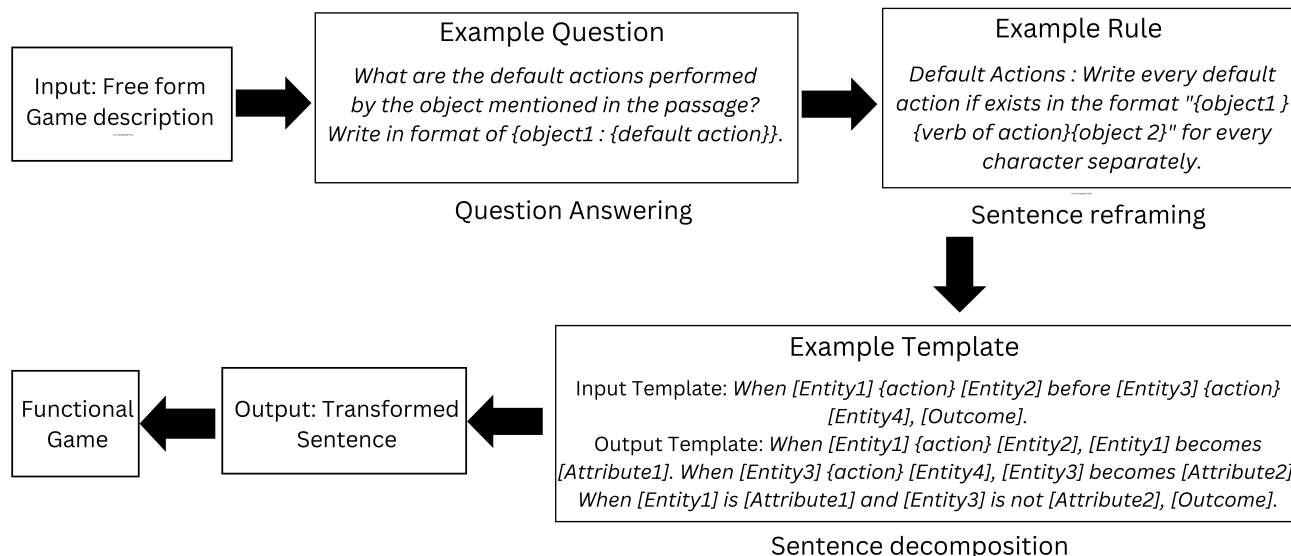


Figure 1. Process flow with example prompt for each step.

A series of iterative tests and comparisons with additional approaches, such as few-shot learning and model fine-tuning, revealed that the suggested strategy performed better overall, especially with unrestricted sentence structures.

Let us consider an input text:
The apricot slows down at border. The rabbit turns into a diamond when hitting a carrot.

Here is a step-by-step trace through the outlined process using the provided input sentence.

- 1) **Question Answering (QA):** The QA component extracts crucial information from the input sentence by asking questions and taking the output in a specific format. It identifies the objects (apricots, rabbits, borders, diamonds), the default actions (apricots and rabbits move), and the conditional actions (speed decrease for apricots, transformation for rabbits).
- 2) **Sentence Re-framing:** Using the information from the above QA, the sentences are then re-framed according to a set of predefined rules that reflect the original free-form sentences. The main goal here is to use a specified set of rules to reconstruct the sentences in a paragraph which are in their basic form in the format subject-verb-object. For example, stating the conditional actions of various objects: when apricots touch a border, their speed decreases, and when rabbits touch a carrot, they turn into diamonds.
Re-framed sentence: If the apricot touches a border, the speed of the apricot decreases. If the rabbit touches a carrot, the rabbit turns into a diamond.
- 3) **Sentence Decomposition:** Next, the Sentence Decomposition step would break down complex sentences into simpler, object-oriented structures. The input would be analyzed to discern patterns of object interactions, such as the apricot’s speed change upon touching a border, and the rabbit’s transformation upon touching a carrot.

An intermediate attribute "mutated" is added while decomposing the sentence resulting in the following sequence of unambiguous sentences [19].

Decomposed sentence (Final Output): If the apricot touches a border, the speed of the apricot decreases. When the rabbit touches a carrot, the rabbit becomes mutated. When the rabbit is mutated, it turns into a diamond.

To sum up our methodology, it offers a comprehensive, structured, and systematic approach to interpret and process natural language text with a high degree of precision and consistency, enabling the user to more accurately describe their intent. Our innovation lies in the strategic application of existing LLM capabilities through a series of system prompts that guide the model to produce outputs in line with specific, predefined rules. This ensures that the transformations maintain the core meaning of the original sentences while stripping away unnecessary complexities, making the text more suitable for generating executable code.

Few-shot learning was initially considered due to its prowess in addressing edge cases with limited data. However, given the vast array of edge cases, rules, and potential issues to address in this domain, few-shot learning proved insufficient. The model would occasionally produce out-of-bound prompts leading to sub-optimal performance. In contrast, our proposed approach, which integrates QA, reframing, and sentence decomposition exhibits robustness against diverse sentence structures, making it an ideal choice for our purpose.

IV. EVALUATION AND DISCUSSION

This section evaluates the performance of the proposed AI-driven assistant in processing 1000 free-form sentences categorized into five types: (1) Grammar/typos, (2) Ambiguous, (3) Unrealizable actions, (4) Overly complex/descriptive, and (5) Non-problematic sentences. Sentences containing grammatical

or typographical errors fall under the first category, "Grammar or Typos" that could cause misinterpretations or inaccurate code translations. The second category, "Ambiguity" refers to statements that have ambiguous references or meanings. Examples of this type of sentence include "It chases it", where pronouns make it difficult to determine exact entities and actions. The third category, "Unrealizable Actions", consists of sentences that describe actions not feasibly translatable into programming logic, exemplified by phrases like "It jumps to heaven". Sentences falling into the "Overly Complex or Descriptive" category are weighed down with too many information or complex structures, which makes it difficult to translate them into concise, executable computer commands. Each of these categories represents a unique facet of the complexity inherent in translating natural language into machine-executable code. The final "Non-problematic sentences" category refers to the sentences which are successfully translatable by the GameChangineer platform into executable code [7] [8] [9]. These sentences are unambiguous and in object oriented structure.

There are several reasons why the final category of "Non-problematic sentences" is included. It serves primarily as a benchmark, providing a point of comparison to assess the efficiency and precision of the AI-powered assistant while processing and interpreting texts that do not present inherent challenges. Furthermore, this category aids in determining whether and how Language Models (LMs) intervention may unintentionally add errors into previously error-free sentences. This will help in evaluating the preservation of sentence integrity after processing and is essential for preserving the overall quality and validity of the research.

The above categorization is based on the platform's algorithms that use symbolic AI to detect grammatical errors, ambiguity, complexity, and unrealizable actions in sentences, indicating potential issues for translating these into executable code. The platform automatically logs the problematic sentences. All logged erroneous sentences are analyzed in this paper.

We discuss the effectiveness of the assistant in identifying and rectifying these issues, thereby enabling accurate translation into executable code. These sentences were written by middle school students with different degrees of experience in both natural language expression and game design when they were first created as parts of game descriptions. This diversity guarantees a wide range of linguistic difficulties, reflective of the intricacies typically seen in natural language programming.

These middle school students received a basic introduction to writing a few simple games with the GameChangineer platform. A small percentage of the students have prior programming experience. However, a vast majority of the students have never programmed before. Participants were given the following instructions to create their game plan: "Write a game plan for creating a game utilizing the available characters."

To ensure the accuracy and feasibility of the translated sentences produced by the LLM, they were given as an input into the GameChangineer platform [7]. This platform provides

a score for each sentence that measures the compatibility with the platform's expected input format [7] [8] [9]. Although some complex sentences can already be decomposed into a sequence of sentences by the GameChangineer platform, it cannot process all the nuances in natural language. We note that all the original problematic sentences were not accepted by the GameChangineer platform.

After the original input sentences were re-written by the LLM using our proposed approach, these new sentences underwent the validation process. Whenever the rewritten sentence(s) are understood with more than 90% certainty by the GameChangineer platform, the conversion will be regarded to have been translated correctly; on the other hand, when it falls below this mark, the output program generated may contain errors. The output program is generated by the GameChangineer Platform. The accuracy and relevance of the LLM-generated results were also assessed manually to ensure the translations effectively communicated the intended meaning. This dual evaluation provides a comprehensive measure of the AI assistant's efficacy in translating complex natural language into machine-executable code by combining automated accuracy assessment with manual semantic verification.

Table 1 presents the results of the sentence categorization from the data-set, highlighting the success rate for each category. The table is divided into three main columns: Sentence Category, Number of Sentences, and Success Rate. These categories include Grammar/Typos, Ambiguous, Unrealizable Actions, Overly Complex/Descriptive, and Non-problematic. Note that when all 5 categories are considered, the success rate was more than 68%. However, if we consider only those first 4 categories (excluding the Non-Problematic category), the success rate by our approach is more than 60%.

The category of Unrealizable Actions encompasses game descriptions featuring actions that are either illogical or infeasible within the game context. This category was the least represented in student written game descriptions, accounting for its minimal proportion. All erroneous sentences are logged by the platform and the distribution of categories reflect the distribution of error types. The presence of actions that cannot be executed within the game complicates the task of rewriting such sentences, contributing to the lowest success rate observed for this category.

A sample of successful and unsuccessful transformations are listed in Table 2. In the following, we will explain a few specific cases.

Let us first examine a successful transformation:
Input Sentence: "Each carrot you eat it gives you 10 speed of

TABLE I. SENTENCE CATEGORIZATION RESULTS

Category	Sentence proportion	Success Rate
Grammar/Typing mistakes	22%	75.6%
Ambiguous Sentences	18%	57.7%
Unrealizable Actions	11.6%	50%
Overly Complex/Descriptive	28.4%	59%
Non-problematic sentences	20%	94%
Total Success Rate	100%	68.4%

TABLE II. INPUT AND OUTPUT SENTENCES

Input Sentence (Free-form)	Output Sentence	Successful?
fox explodes rabbit when it is the rabbit is touched by it.	If the rabbit is touched by the fox, the fox explodes.	Successful
When the fox touches the carrot before the rabbit, the fox dies.	When a fox touches a carrot, the fox becomes happy. When a rabbit touch the carrot, the rabbit becomes happy. When the fox is happy and the rabbit is not happy, the fox dies	Successful
if the fox sees the rabbit and the fox chases it and catches it the rabbit explodes.	When the fox sees the rabbit, it chases the rabbit. When the fox catches the rabbit, the rabbit explodes.	Successful
Rabbit move with wsad.	When the W key is pressed, the rabbit moves up. When the S key is pressed, the rabbit moves down. When the A key is pressed, the rabbit moves left. When the D key is pressed, the rabbit moves right.	Successful
The kitten must touches all the cakes to win without having the fox touch the kitten.	When the kitten touches all the cakes, the player wins the game. If the fox touches the kitten, the player loses the game.	Successful
When cake is touched by rabbit the speed of the rabbit is 0.5 for 3 seconds and cake disappears.	When the rabbit touches the cake, the rabbit's speed becomes 0.5 for 3 seconds. When the rabbit touches the cake, the cake disappears.	Successful
When the ball touches the bottom border it gains 1 point.	When the ball touches the bottom border, the score adds 1.	Successful
the pointer shoot the ball with the spacebar.	The pointer shoots the ball when the spacebar is pressed.	Successful
The rabbit dies when it touches the diamond for 5 seconds.	When the rabbit touches the diamond for 5 seconds, the rabbit becomes dead. When the rabbit is dead, it dies.	Successful
The rabbit must flee from the fox.	If the rabbit sees the fox, the rabbit flees.	Successful
the sheep is trying to eat all the broccoli.	If the sheep sees broccoli, the sheep eats broccoli.	Successful
When a tiger is shot for the 2 times, it disappears.	If a tiger is shot twice, it becomes mutated. When a tiger is mutated, it disappears.	Successful
if you do not touch the apricot within 3 minutes game is over.	If the apricot is not touched within 3 minutes, the game is over.	Unsuccessful
If the donut is blue and W is pressed, the donut moves up	There is a donut. When the W key is pressed, the donut moves up 10. If the donut is blue, then the color of the donut changes.	Unsuccessful
fox explode 20 seconds after start.	There is a fox. If the start is 20 seconds, the fox explodes.	Unsuccessful
The panda kills foxes when it reaches a score of greater than or equal to 14.	If the score is greater than or equal to 14, the score adds 1.	Unsuccessful
There are 3 snakes going side to side at 3 pixels per frame.	The snakes move side to side at 3 pixels per frame.	Unsuccessful

the rabbit."

Transformed Sentence: "When the player eats a carrot, the speed of the rabbit increases by 10."

The rewritten sentence establishes a clear subject "the player" and action "eats a carrot" by explicitly stating "When the player eats a carrot" as the condition. The usage of "it" in the original sentence lacks a clear antecedent, which could cause parsing issues. This ambiguity is eliminated in the rewritten sentence by defining the action and its result clearly without the use of potentially confusing pronouns. The phrase "the speed of the rabbit increases by 10" explicitly describes the effect of the action in a structured way that is easy for a machine to parse. This approach makes it easy to transform into a computational operation because it explicitly states the change ("increases by 10") and the attribute being impacted ("the speed of the rabbit").

Let us look at a few examples that show how the AI-driven assistant can recognize and fix typos and grammatical problems in sentences.

User Input: "When Dis pressed the bird moves right."

Transformed Sentence: "When the D key is pressed, the bird moves right."

The original sentence contained a typo ("Dis" instead of "D") and lacked a verb, making it unclear. The AI model successfully recognized the typo, corrected it to "D", and inserted the necessary word "key" to complete the sentence.

User Input: "When the rabbit is silver, the speed of the rabbit."

Transformed Sentence: "If the rabbit is silver, the speed of the rabbit increases."

The user input was incomplete, with no clear action described for the rabbit. The AI assistant inferred the most logical action "speed increase" based on the context, making the instruction actionable for programming.

Consider one unsuccessful transformation which could not adhere to the intent of the user.

User Input: "When the rabbit is yellow, the rabbit touches the fox, the fox explodes."

Transformed Sentences: "If the rabbit is yellow, the rabbit touches the fox. If the rabbit touches the fox, the fox explodes."

Intended Output: " If the rabbit is yellow and the rabbit touches the fox, the fox explodes."

The input is incorrectly divided into two conditional statements by the original transformation. This method falsely implies that the fox's explosion is a two-step process that depends on the requirements being satisfied sequentially as opposed to simultaneously. The intended output, on the other hand, combines the two conditions into a single compound condition meaning that the fox will explode if both conditions are met simultaneously and directly. This showed that the input sentence is ambiguous and the AI-assistant could not successfully transform the sentence.

Let us look at an unsuccessful example in the fifth category, Non-problematic sentences.

User Input: "When a ball sees the rock, the ball flees from the rock."

Transformed sentences: "When the ball sees the rock, the ball becomes scared and flees from the rock."

The transformed sentence is considered unsuccessful here,

primarily due to the addition of an unwanted attribute "scared" to the output sentence. This is an example where the LLM hallucinated leading to add an extra and unnecessary attribute [20]. Such hallucinations can significantly impact the utility and accuracy of LLMs, especially in applications requiring strict adherence to input data without the addition of interpretative or speculative elements. LLMs occasionally "hallucinate," or provide missing information [20]. We found that unsuccessful conversions due to hallucination account for 6% of Non-problematic sentences. For the problematic sentences in the other four categories, hallucination is responsible for about 12% of the unsuccessful transformations.

We did not compare our results with LLM based code generation platforms such as Copilot [21] because our goal is to rewrite erroneous sentences so that they become clear and unambiguous. On the other hand, while Copilot may be able to generate code on an erroneous sentence, it generates the code by its own interpretation arbitrarily. In addition, GameChangineer can process hundreds of sentences at a time, but the user must interface Copilot differently by feeding a few sentences at a time.

V. CONCLUSION AND FUTURE WORK

This paper presents a method of converting free-form natural language sentences into a sequence of unambiguous, simplified sentences that can subsequently be translated into machine-executable code. The utilization of LLMs has shown promise in addressing the inherent difficulties brought about by verbosity, ambiguities, complexity, and possible errors. Our approach, which combines aspects of Question Answering, Sentence Reframing, and Sentence Decomposition has demonstrated a notable capacity to handle a wide variety of linguistic patterns and semantic complexities. More than 68% of the 1000 problematic and non-problematic sentences were correctly converted by the proposed method.

There are areas for improvement, particularly in understanding complex conditional relationships and refining the LLM methodologies, aiming to reduce the incidence of hallucinations. The results highlight the inherent challenges of processing natural language, particularly in dealing with the nuances of human language. Additionally, they draw attention to how AI-powered systems have the potential to greatly enhance our comprehension and interpretation of words with unclear structures, which is an important area of study in the field of natural language programming.

ACKNOWLEDGEMENT

This research is supported in part by NSF grant 2101021.

REFERENCES

- [1] C. Yang, Y. Liu, and C. Yin, "Recent advances in intelligent source code generation: A survey on natural language based studies," *Entropy*, vol. 23, no. 9, p. 1174, 2021.
- [2] D. Baidoo-Anu and L. Owusu Ansah, "Education in the era of generative artificial intelligence (AI): Understanding the potential benefits of chatgpt in promoting teaching and learning," *SSRN Electronic Journal*, January 2023, published by Elsevier BV, retrieved: March, 2024, Available at SSRN: <https://ssrn.com/abstract=4337484> or <http://dx.doi.org/10.2139/ssrn.4337484>.
- [3] T. P. Tate, S. Doroudi, D. Ritchie, Y. Xu, and M. W. Uci, "Educational research and AI-generated writing: Confronting the coming tsunami," January 2023, published by Center for Open Science, retrieved: April, 2024. [Online]. Available: <https://doi.org/10.35542/osf.io/4mec3>
- [4] D. Mogil et al., "Generating diverse code explanations using the GPT-3 large language model," in *ICER '22: Proceedings of the 2022 ACM Conference on International Computing Education Research*. Association for Computing Machinery, 08 2022, pp. 37–39.
- [5] F. F. Xu, B. Vasilescu, and G. Neubig, "In-IDE code generation from natural language: Promise and challenges," *ACM Trans. Softw. Eng. Methodol.*, vol. 31, no. 2, mar 2022, retrieved: April, 2024. [Online]. Available: <https://doi.org/10.1145/3487569>
- [6] C. Niklaus, "From complex sentences to a formal semantic representation using syntactic text simplification and open information extraction," Ph.D. dissertation, 03 2022, retrieved: April, 2024. [Online]. Available: <https://opus4.kobv.de/opus4-unipassau/frontdoor/index/index/docId/1054>
- [7] M. S. Hsiao, "Automated program synthesis from object-oriented natural language for computer games," in *Proceedings of the Controlled Natural Language Conference*, August 2018.
- [8] —, "Multi-phase context vectors for generating feedback for natural-language based programming," in *Controlled Natural Language*, September 2021.
- [9] —, "Automated program synthesis from natural language for domain specific computing applications," Patent 10 843 080, November, 2020.
- [10] F. Shi, D. Fried, M. Ghazvininejad, L. Zettlemoyer, and S. I. Wang, "Natural language to code translation with execution," 2022, arXiv.
- [11] T. H. M. Le, H. Chen, and M. A. Babar, "Deep learning for source code modeling and generation: Models, applications, and challenges," *ACM Computing Surveys*, vol. 53, no. 3, p. 1–38, Jun. 2020, retrieved: April, 2024. [Online]. Available: <http://dx.doi.org/10.1145/3383458>
- [12] A. Vaswani et al., "Attention is all you need," in *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2017, pp. 5998–6008, retrieved: April, 2024. [Online]. Available: <https://arxiv.org/abs/1706.03762>
- [13] M. L. Zong and B. Krishnamachari, "A survey on GPT-3," 2022, arXiv, retrieved: April, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:254221221>
- [14] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *North American Chapter of the Association for Computational Linguistics*, 2019, retrieved: April, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:52967399>
- [15] Y. Liu et al., "Roberta: A robustly optimized BERT pretraining approach," 2019, CoRR, abs/1907.11692, retrieved: April, 2024. [Online]. Available: <http://arxiv.org/abs/1907.11692>
- [16] J. Wei et al., "Chain-of-thought prompting elicits reasoning in large language models," in *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35. Curran Associates, Inc., 2022, pp. 24 824–24 837.
- [17] J. Akanya and C. G. Omachonu, "Meaning and semantic roles of words in context," *International Journal of English Language and Linguistics Research (IJELLR)*, vol. 7, pp. 1–9, 03 2019.
- [18] M. Rosol, J. S. Gasiior, J. Łaba, K. Korzeniewski, and M. Młyńczak, "Evaluation of the performance of GPT-3.5 and GPT-4 on the polish medical final examination," *Scientific Reports*, vol. 13, no. 1, p. 20512, 2023, retrieved: April, 2024. [Online]. Available: <https://doi.org/10.1038/s41598-023-46995-z>
- [19] N. Wies, Y. Levine, and A. Shashua, "Sub-task decomposition enables learning in sequence to sequence tasks," 2023, arXiv.
- [20] Z. Ji et al., "Survey of hallucination in natural language generation," *ACM Computing Surveys*, vol. 55, no. 12, p. 1–38, Mar. 2023, retrieved: April, 2024. [Online]. Available: <http://dx.doi.org/10.1145/3571730>
- [21] GitHub, "About github copilot," 2024, retrieved: April, 2024. [Online]. Available: <https://docs.github.com/en/copilot/about-github-copilot>