# Code Generation Empowered by Natural Language Processing and Machine Learning Algorithms

**Dr.S.Girirajan[1*], S.Vidhya[2], S. Tamilselvi[3], S.Durai[4], R. Senthilkumar[5], Dr. S. Vinoth Kumar[6]**

[1]Assistant Professor, Department of Computing Technologies, School of Computing,
SRM Institute of Science and Technology, Kattankulathur, Chengalpattu-603203, Tamil Nadu, India. Corresponding author Email: girirajans.cse@gmail.com

[2]Assistant Professor, Dept of Information Science and Engineering, The Oxford College of Engineering, Bangaluru, India. Email: vidhyatoce@gmail.com

[3]Assistant Professor, Department of Artificial Intelligence and Data Science, Panimalar Engineering College, Chennai, Tamilnadu. Email: stamilselvi@panimalar.ac.in

[4]Associate Professor, Department of Computer Science and Engineering, Vel Tech Rangarajan Dr.Sagunthala R&D Institute of Science and Technology, Avadi, Chennai, Tamilnadu-600062. Email: duraitrichy@gmail.com

[5]Senior Assistant Professor, Computer Science and Engineering, Alva's Institute of Engineering and Technology, Karnataka. Email: senthil2481@gmail.com

[6]Professor, Department of CSE, Vel Tech Rangarajan Dr.Sagunthala R&D Institute of Science and Technology, Chennai, Tamilnadu. Email: profsvinoth@gmail.com

**Abstract:**

The goal of this study is to revolutionize code creation processes by investigating the synergistic union of machine learning (ML) and natural language processing (NLP). Enterprising non-programmers with entrance barriers, traditional approaches to code generation frequently demand expert-level programming expertise. Development teams can communicate coding tasks in natural language by utilizing NLP techniques like language modeling and semantic parsing. This helps to close the gap between human intent and instructions that can be executed by a computer. By incorporating ML techniques, the system may also more effectively understand and produce code that is compatible with a wider range of programming languages and paradigms. This research clarifies the revolutionary potential of NLP and ML-driven code creation and highlights its consequences for software development efficiency, accessibility, and innovation through an extensive assessment of current developments and case examples.

**Keywords**: Code creation processes, Machine learning (ML), Natural language processing (NLP), Programming expertise, Software development innovation.

## 1.    INTRODUCTION:

In the software industry, human beings write the majority of the project documents. One of the crucial phases in a software lifespan model is this procedure. The improved document will modernize procedures and save project expenses for upkeep. Software undertaking Metadata preserves specifics of the application development cycle's requirements, evaluation, programming, testing, and installation processes. The sole focus of this research project is on automatically creating programming documentation. Code instructions for the procedure for maintenance should be concise, easily understood, and devoid of any confusion for engineers. Documentation typically requires a significant amount of human labor and effort. The quantity of pages in documentation per

kilogram one statistic used to estimate the cost of a project involving software is code. A key component of the software initiative's edition upgrade will be successful documentation.

This suggested automated source code reference system uses natural language processing (NLP) techniques to cut down on the period and expense of documentation development. Unlike other languages that use object-oriented programming like Java, C has a simple grammar set, which allows this suggested system to produce a reference for C programs. Using a straightforward syntax-based C-lite programming language and compilation, this research project creates an initial design for a computerized code description process. Few code-based tools currently in existence, such as JavaDoc, which only creates a reference for Java's standard techniques by pulling data from the official documentation for the language itself, are limited in their functionality. Using Natural Learning Processing techniques, the system shown in this work can automatically produce documentation for C programs. [1] The enhanced technique contact graph methodology is used by this system to automatically generate two main forms of writing: (i) method-based records, which documents defined by users and predetermined methods. This allows for the explanation of defined by users' methods via Natural Language Generation grammatical structures, while established methods can be recorded through the C programming language compiler documentation (ii). Statement-based documentation documents each program line, including any methods that may be present in the program in question.

The suggested system generates a description as a result by using a C program as a source of input. Program complexity varies depending on whether user-defined or established techniques are present. Free of context Languages are constructed for the C programming the language's semantics. Program statements including gestures, illnesses statements, method terms, and so forth can be recognized by this set of criteria. Next, utilizing the NLG approach, every statement type is transformed into a line specification with the required identifiers. This system's Source Word Usage Model, which serves as the central component of the application code description process, was created using a unique method.

### 1.1. The limitations of traditional source code methods for modelling:

Explained four conventional methods for addressing the lexical and psychological aspects of origin code:

(a) Simple neural language models;

(b) Probabilistic grammars;

(c) domain-specific language-guided models;

(d) Simple stochastic models of languages.

These conventional methods still have several significant problems, which DL models can successfully address. [2] This section goes over each strategy in more depth, including its drawbacks.

Domain-specific Language Guided Models: The parameters criteria and assertions that are used to create a program's structure are typically defined using Domain-Specific Languages (DSLs). For typical code statements (such as controlling movement, feedback, and parentheses), DSL-based models generate distinct grammatical rules. A DSL's syntax is typically smaller than that of a language used for all purposes, making it more effective for certain code-producing jobs. By employing a calculus of concise kinds, the DSL-based approach narrowed down the searching area for Java expression recommendations and created a more complex functional script fulfillment.

Grammars with Probabilities: Generation rules in official languages specify the number of strings that can exist. Free-of-context grammatical structures are a collection of manufacturing rules that are

applicable in all circumstances; that is, a manufacturing rule having only one non-terminal sign on the opposite side. A common method for defining a programming language's architecture is CFG, which may be used to create an Abstraction Grammar Tree from the original code.

**Models of n-gram Languages:** In addition to the two grammatical models mentioned above, n-gram is a straightforward yet powerful empirical communication model. To be more precise, this framework claims that each symbol or phrase is critically reliant on the n minus one character or phrases that came before it, as shown by the equation that follows: It is easy to calculate, WHERE by just adding up every instance of each n-gram in the initial set. As the relationships and rules of the language of programming are readily absorbed from the code itself, n-grams have a clear generality over other structural models (e.g., stochastic grammatical structures and DSL-guided modeling).

Basic Models of Neural Programs:

One hidden layer neural network embed models first translate a word's immediate encoding into a second word-embedding vector that is considerably less in length (e.g., 101–1010) than its vocabulary size, as opposed to directly adding the average number of each preceding token. Spread representation of words is another term for this concept. A fully integrated neural network with a single concealed layer was fed an array of up to n - 1 prior tokens about the current word in its original work. The likelihood of the following word was determined at the result layer using the function SoftMax.
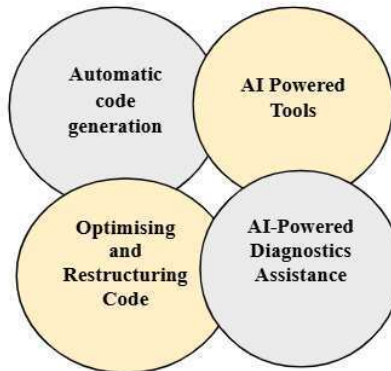
## 1.2.    Programming's Development via Integration of Machine Learning (ML):

The process of creating, implementing, and maintaining software has significantly changed as a result of the incorporation of machine learning (ML) into coding. Machine learning models were initially created and trained independently of the software that they were eventually linked to. Using tools like sci-kit-learn or tensor flow developers would build machine-learning models and then incorporate them into their apps. This method frequently needed a large amount of engineering work and required manual integration. Libraries appeared as machine learning (ML) spread throughout different applications, enabling developers to include ML features straight into their programs. Arrangements such as Tensor Streaming and TensorFlow Lite were developed for mobile platforms making it simpler to include machine learning algorithms into software programmes. Model selection, parameter adjustment, and even designing features may now be automated by engineers with the emergence of AutoML tools. This simplified the process of integrating the ability to learn into apps for software engineers without a lot of experience with ML. More fluid and flexible systems with learning and evolution capabilities are probably in store for programming's use of machine learning in the future. Software applications may automatically adjust to shifting circumstances and needs without the need for human intervention because of continuous learning approaches, which enable algorithms to gradually update and enhance it as additional information becomes accessible.

## 1.3.    Programming using Machine Learning Applications:

The process of employing machine learning and artificial intelligence methods to instinctively produce original code or code fragments based on the highest level demands, requirements, or samples is known as "automated code generation." This method seeks to increase efficiency and speed up the creation of software by lowering the amount of programming by hand required. Artificial intelligence (AI) models that have been educated on enormous volumes of written information, like OpenAI's Codex, can comprehend basic explanations of computing duties or needs. Pre-defined blueprints, patterns of design, and code fragments are frequently used by automatic code generation technologies to produce code. The code produced will be organized and adhere to company norms due to these templates, which capture conventional programming constructs and

best practices. Some automatic code generation systems use machine learning models to discover connections, trends, and meanings from current code that have been taught in big code libraries. These models can then generate new code that resembles the style and logic of the training data. AI-powered code generation can be tailored to specific domains or use cases, such as web development, mobile app development, data analysis, or machine learning. By understanding the context and requirements of a particular domain, automated code generation tools can produce more relevant and efficient code. Code Optimization and refactoring are essential practices in software development aimed at improving the performance, readability, maintainability, and overall quality of code.



**Figure 1.1:** Machine Learning Applications in Programming

These technologies use machine learning algorithms that have been built on massive archives of code to identify trends and rewrite code appropriately. They can, for instance, spot redundant code, remove classes or functions that are used repeatedly, and recommend improved labeling of variables or standards for coding. Artificial intelligence is often used to autonomously produce efficient code. Artificial intelligence (AI) methods, for instance, can offer outstanding performance software that is suited to specific computer designs or standardization goals in the field of numerical computation or neural networks. This covers methods like code concentration, auto-vectorization, and circuit rolling. Artificial intelligence (AI) technologies are becoming more and more important in streamlining and enhancing the crucial tasks of bug discovery and repair. In the combination of artificial intelligence and software development, machine learning (ML) offers a range of apps, as Figure 1.1 illustrates:
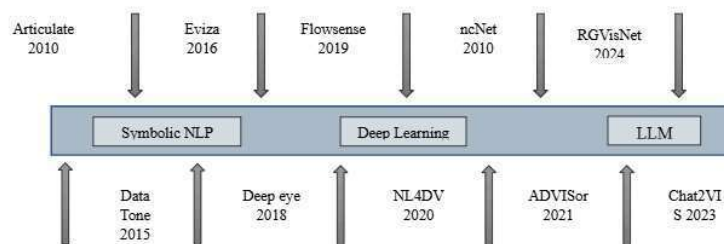
Through the study of the source code without operation, artificial intelligence-powered static analysis instruments can discern potential mistakes, weaknesses, and code odors. Using artificial intelligence methods taught on large code databases, these tools identify trends that may indicate frequent mistakes in programming, vulnerabilities in security, or speed issues. AI-powered systems may generate information that is useful, constant, and understandable by programmers with various degrees of skill by utilizing Natural Language Generation (NLG). Documentation specific to a given repository or coding activity can be produced by AI techniques. Development efficiency is increased with contextual documentation, which also lessens the mental strain that comes with exploring complicated frameworks. Script searches driven by AI can offer condensed summaries of code fragments that are obtained from discussion boards or libraries.

## 2.    RELATED WORKS:

Several superior code generation models have surfaced lately as a result of ongoing improvements to the universal code generation technology that relies on natural language description. Still, the assessment of the standard of the created code is relatively new. Code quality can be assessed in two ways. The first involves contrasting the code to the standard code to see how comparable it is. [4] Its

quality increases with similarity level. Examining if the program performs the features specified in the specifications is another way to gauge the standard of the created code. Recent research has typically employed the initial technique, which involves calculating the degree of resemblance between the created language and the standard source. The sole measure of assessment they employ is the token level's BLEU. The sequence level resemblance measure strACC that CNN suggests didn't make a difference. These approaches merely analyze its code as plain written content, disregarding the architectural aspects of the code as computer languages.

The natural language processing computational methods that form the foundation of NL2VIS can be categorized broadly into two groups: the traditional symbolic-based NLP approaches, which rely on illustrations of language framework and clear guidelines, and the newly developed neural network translation methods, which rely on linguistic models built primarily via neural networks. The development of NL2VIS algorithms in printed works is shown in Fig. 2.1, [5] wherein the authors argue that ultimately making use of the most advanced AI systems, such as master's degrees, for comprehension of languages and creating code is the logical path. The development of advanced machine-learning approaches is also displayed.



**Figure 2.1:** Recent changes in NL2VIS systems and the proposal for integrating Chat2VIS in a timeline format.

The creation of certain linguistic software comments is the focus of software remark creation. Converting the language of programming code into verbal remarks expressed in normal tongues is the main job. In the meantime, notes can explain both the features and the programmer's architectural intentions found in the actual code. [6] In a nutshell, computerized code remarking is the process of automatically generating written descriptions for source code using linguistic analysis. This process can disclose information about program reasoning, architectural intentions, performance, and the significance of associated variables, among other things.

Using a machine learning model, the creators of Deep Coder: Training to Write Programmes have built a program generation machine that produces the programs. It explains how the model for machine learning works, and how it is used for program generation to reduce estimations. The researchers demonstrated how using machine learning approaches for algorithm creation greatly decreases the number of stages needed to build the software and speeds up the time required for the execution of software manufacturing. Synquid is yet another instance of a contemporary program generation system. [7] The software synthesizer based on improvement classes is called Synquid. Synquid uses iteration kinds, which are specified, to produce functionally compatible code. Cognitive limitations, structural limitations, and search methods are the three main elements of software creativity. The way a consumer engages with technology is defined by a behavioral restriction. The program's exploration space is restricted by physical constraints.

The natural language components are executed predictably as formally as what they refer to as "realistic computing." The setting, confusion, communication, words, anaphoric interactions, deixis, and others are a few examples of characteristics of natural speech. [8] Automated code generation, set aside phrases, established syntax, dictionaries of data, established software pieces, support for

various languages of programming, passive citations, business-, acquiring-, and documentation-focused features were the features that they presented. Subtle allusions, English-like communication, databases of data, focused on learning, code development, and reporting production were the subsequent ways they collected languages.
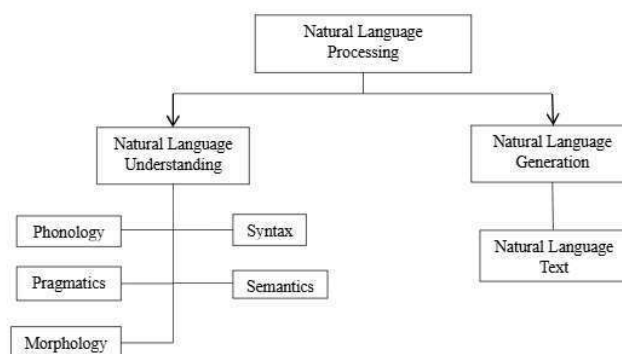
Source Code Synopsis:

Relevant coding for presenting work is generally categorized as template-driven or data-driven. The data-driven approach that is the most directly associated with this study is recent work. The proposal in that study is to label the words in the original code using the AST, and then input the annotated version into a sequence-to-sequence artificial neural model. The primary innovation lies in the AST-annotated form known as Substance-based Moving; the model is a ready-made encoder-decoder. [9] SBT is a method for making sure that every word in the source text is related to the sort of AST node they belong to and for smoothing the AST. As an illustration, the code request. remove(id). The idea is for the terms "request," "remove," and "id" to be understood about the sentences in which they occur. Here, a node for Procedure Invoked. This method depiction serves as an essential point of reference for reference in the studies that follow. Although SBT was demonstrated in that study to achieve an impressive thirty-eight BLEU efficiency, we warn that this isn't identical to the outcomes of our trials.

## 3. METHODS AND MATERIALS:
### 3.1. NLP:

NLP is divided into two categories: the generation of natural language and comprehension, which progresses towards the goal of creating texts and knowledge. The general NLP description is shown in Figure 3.1. The discussion of the creation of natural languages (NLG) and understanding (NLU) is the aim of this chapter. Supporting one or more methods' or sets' specializations is the aim of natural language processing (NLP). It is possible to incorporate creating languages along with comprehension using the NLP assessment metric on an algorithmic system. Linguistic event detection even makes use of it. It is designed to be an original flexible system with a distinct pipeline for each language that is used to extract events from texts in the language, and Italian. A versatile collection of top-notch international natural language processing tools is integrated into the framework. More complex tasks like cross-lingual identified object connection, semantics role labeling, and time standardization are all integrated into the pipeline along with components for fundamental NLP processing.

As a result, the cross-lingual structure permits the understanding of duration, locations, individuals, and occurrences in addition to the relationships among them. The purpose of each pipeline's output is to serve as a starting point for an algorithm that generates event-centric information graphs. Every module takes in standard input, annotates it, and outputs standard output, which is then used as the input for pipelines to the subsequent modules. Because their workflows are designed with a data-driven construction, modules can be easily added or removed. The modular period design additionally permits various combinations and dynamic distribution.

**Figure 3.1:** NLP Classification

### 3.2. NLU:

NLU makes it possible for computers to comprehend natural language and analyze it by removing ideas, organizations, feelings phrases, and other elements. In customer service applications, it's utilized to comprehend issues that consumers report, in person or text. The study of language significance, linguistic setting, and variations in language is known as semantics. Therefore, it is critical to comprehend the various NLP levels and key terms. Next, we go over a few terms that are frequently used at various NLP levels. Linguistics phonology is the branch that deals with the orderly organization of sounds. The word "phonology" originates in Ancient Greek, where "phono" denotes sound or noise and "-logy" denotes speech or language. According to More (1993), phonology is "the study of sound and everything related to the framework of language," but it also widely refers to the sounds that make up language and deals with the behavior and structure of sounds as a sub-discipline of semantics. The linguistic utilization of sound to store meaning in any human language is known as phonology. The word's constituent parts stand in for tiny units of measurement, or morphemes in the morphemes are the starting point for shape, which makes up the structure of words. The term pre-cancellation, for instance, can be morphologically examined to reveal three distinct phases: the beginning before, the initial cancella, and the ending -tion. This is an example of a morpheme. [10] All words have the same meaning when it comes to phrases; humans may break down any word that is unfamiliar into its constituent phrases to comprehend its meaning. When a verb is suffixed with -ed, for instance, it indicates that the operation of the verb occurred in the past. Lexical forms are words (such as tables and chairs) that are not separable and have independent meanings.

Both humans and NLP algorithms decipher word meanings in lexical systems. Phrase-level comprehension is granted via various process kinds; one of which involves assigning a part-of-speech label to every word. Words that can possess as many parts of speech are given the most likely part-of-speech tag in the above process depending on the setting in which they appear. Single-meaning words can take the place of semantic models at the level of the lexicon. In actuality, the type of representations in the natural language processing system differs depending on the semantics concept that is used. As a result, the word structure is analyzed at the level of vocabulary about its lexical significance and PoS.

Words are classified into expressions, sentences to provisions, and ideas to sentences at the grammatical level following PoS tagging at the level of lexicon. By dissecting the language's grammatical framework, it highlights how a sentence should be formed correctly. This level produces a sentence that displays word structure relationships. It also goes by the name of processing, which uncovers phrases that have more significance than a single word. The grammatical level looks at things that the level of language ignores, like word arrangement, stop-words, shape,

and the PoS protocol. Word dependency will alter as words are arranged differently, and this may also impact how well sentences are understood.

Finding a sentence's correct meaning is the most crucial effort from a semantic perspective. Humans rely on their understanding of syntax and the ideas contained in a sentence to determine its meaning; machines are not dependent on these abilities. By analyzing the logical structure of an expression and selecting the most appropriate phrases to comprehend how various ideas or phrases interact with one another, semantic processing ascertains the potential meanings of a given sentence. The semantic level looks at a word's dictionary definition or perception based on the sentence's context.

The discourse stage of NLP works with a maximum of a single phrase, whereas the semantic and syntax layers deal with paragraph-length entities. By establishing connections between phrases and words that ensure coherence, it analyses conceptual structure. The functional level is concerned with information or material that doesn't exist in the written work itself. It addresses inferences made by viewers and by the speaker's words. The phrases that are not spoken directly are examined by it. Understanding the topics covered in the text requires an understanding of the real world.

### 3.3.    Methodologies based on semantic parsing:

Semantic parsing is a branch of the processing of natural languages that converts natural language speech into forms that are comprehensible to machines. The ultimate goal is to derive exact meaning so that robots can react to these speech outputs. It is used in many different fields, such as automated reasoning, taxonomy induction, translation by machine, query responding, and creating codes. In a neural semantic analysis for code generation, the Context-Based and Parameter Cloning approach incorporates context-specific data to improve identification. This method makes use of the software context when decoding by using an encoder-decoder device. A supervised copy method duplicates surrounding signs, even if they were not seen throughout training, while an additional attention technique aligns language terms with surrounding markers. Language structure has an important role in shaping the decoding procedure of neuronal semantically processed models. Syntax restrictions are crucial, as demonstrated by methods such as imposing type restrictions in question creation and generating information database queries sequence-based. [11] They introduced a method for generating code from process to sequence that makes use of a pointer and computational models to copy phrases from inputs. The suggestion is to use abstract trees of syntax to make programming languages for general use easier to understand.

Symbolic reasoning and recognition of pattern modeling are two further methods for applying semantic parsing techniques. In program synthesis, the sketch is the process of using unfinished programs or designs to articulate high-degree descriptions. Their methodology centered on a framework for trust modeling that rigorously defines various types of ambiguity, such as input, data, and model concerns. A score of trust is produced by including these measurements of certainty as variables in a model of regression. This score improves the model's overall comprehension by helping to pinpoint the sources of forecast uncertainty. This approach is significant because it establishes the foundation for interactive computing techniques, which enable programs to involve users in enlightening conversations when faced with ambiguous or incomplete data. This method works especially well for fixing mistakes or insufficient user inquiries that arise after processing. Conversational artificial intelligence plays a crucial role in filling in the gaps and flaws in logical forms by facilitating an interactive dialogue between the user and the program. Iterative swaps not only fix errors but also fill in lacking details. Furthermore, such dialogue plays a major role in reducing the size of the search field, which leads to results forecasts that are more accurate and trustworthy.

### 3.4.  Module for Code Summarization:

Firstly, we describe time-related interactions among words using a Simultaneous Short-Term Long-Term System-based recurrent neural network based on the original code embedding sequences. There are three exponential gates in a linear sequential transistor cell. It creates a new concealed state at each time step t by combining the last concealed state with the encapsulation in the phrase in the initial code sequence. The LSTM cell's gates determine each input's rate of transmission. The concealed information changes as follows at every increment t:

$$j = \alpha \ (V_j g_{t-1}^f + W_i F_y(xt) + a_i \tag{1}$$

$$e = \alpha \ (V_e g_{t-1}^f + W_e F_y(xt) + a_e \tag{2}$$

$$p = \alpha \ (V_p g_{t-1}^f + W_P F_y(xt) + a_p) \tag{3}$$

$$s = \sin h \ (V_s \ g_{t-1}^f + W_s F_y(xt) + a_s) \tag{4}$$

$$s_t = e \odot s_{t-1} + j \odot s \tag{5}$$

$$g_t = p \odot \sin h \ (st) \tag{6}$$

Where are weighting matrices for the last concealed stages;    value matrices for the original code embedded data;    are inclinations; and σ represents the element-wise sigmoid value and ⊙is the element-wise summation. To keep things simple, we refer to the computation above as follows.

$$g_t^f = LSTM_y \ (g_{t-1}^f) \ F_y(xt) \tag{7}$$

Where represents the step's hidden state in the Bi-LSTM encoder for x. As only gets data from code before the position t, we additionally integrate data following location t into concealed states using Simultaneous LSTM:

$$\overrightarrow{g_t^f} = \overrightarrow{LSTM_y} \ (\overrightarrow{g_{t-1}^f}) \ F_y(xt) \tag{8}$$

$$\overleftarrow{g_t^f} = \overleftarrow{LSTM_y} \ (\overleftarrow{g_{t-1}^f}) \ F_y(xt) \tag{9}$$

$$g_t^f = \left[\overrightarrow{g_t^f}; \overleftarrow{g_t^f}\right] \tag{10}$$

Next, we utilise an additional LSTM-based parser to produce a written brief:

$$g_t^c = LSTM_y \ (g_{t-1}^c, \ [F_x(x_{t-1}); d_{t-1}^c]) \tag{11}$$

Where $g^f, g_0^c = \left[\overrightarrow{g_n^f}; \overleftarrow{g_f^l}\right], (x_{t-1})$ the produced phrase of the $t - 1^{th}$phase marks the combination functioning and $g_t^c$ is the secret state from the $t^{th}$ n LSTM parser where $d_{t-1}^c$, is the current context vector generated through the conventional attentiveness system:

$$\sigma_{jt} = \frac{EXP \ (e(g_j^f, g_t^c))}{\sum_{j-1}^n EXP \ (e \ (g_j^f, g_t^c))} \tag{12}$$

$$d_t = \sum_{j=1}^n \sigma_{jt} \ g_j^f \tag{13}$$

$$w_t^f = V_w \ [g_t^c, d_t] + a_w \tag{14}$$

To determine the relationship between concealed states, e is an adaptable function. An adaptable matrix is used, and the function we capitalize on is a basic bi-linear polynomial. [12] Lastly, the linear layer receives the background vector and parser output combined, resulting in the resultant word dispersion $R_\gamma$:

### 3.5. Module for Code Creation:

To improve the performance of our primary task, code summarization, we build an additional task with an opposing network able to produce an original code series from its summaries. Since the entire structure is essentially identical to the source code for the synthesis module, we will just refer to it by this name:

$$\hat{g}_t^f = BiLSTM\ \hat{g}_{t-1,}^f (F_x(xt)) \tag{1}$$

$$\hat{g}_t^c = LSTM_y(\hat{g}_{t-1,}^c\ [F_x(x_{t-1}); \hat{d}_{t-1}^c]) \tag{2}$$

Where the terms $g_t^f$ and $g_t^c$ refer to the secret states that comprise the $t^{th}$ phase in the LSTM detector and the Bi-LSTM encoder, respectively; whereas $g_t^c$ represents the setting vector, x't stands for the created phrase of the $t^{th}$ step.

The code creation generator and the program summation receiver both attempt to operate with the original code input, therefore they think that dividing up the LSTM-cell parameters amongst them would simplify the model. For the code summation coder and the code generation decoder, this is also true. It is observed that there is minimal difference in performance between the use of separate encoding devices and decoding devices.

Additionally, the reduction function is intended to equal the target word, x't negative logarithmic chances:

$$\log_{ds}(y, x, \emptyset_{ds}) = -\sum \log\ (\hat{R}_\gamma\ (xt)) \tag{3}$$

Where all of the parameters that can be trained in the code generating module are indicated by $\emptyset_{ds}$ and $d_{t-1}^c$, and $R_\gamma$ is the produced word dispersion determined by $g_{t-1,}^c$

## 4. RESULT ANALYSIS AND DISCUSSION:
### 4.1. Dataset:

**Table 4.1:** A synopsis of the training set

| Languages used for Programming | Total Number of Files | Total Number of Lines | Count of Repositories |
|---|---|---|---|
| C | 1378 | 200 | 4899 |
| Python | 1341 | 267 | 58629 |
| Java Script | 1875 | 707 | 23762 |
| Type Script | 467 | 85 | 3863 |

We gather a sizable source code collection to train and assess the sequence of the code conclusion model. It includes more than 1.1 billion lines of origin code in the programming language Python, C#, the case of JavaScript, and Typescript programming languages, as listed in Table 4.1.

A selection of more than 51500 top-starred initiatives on GitHub has been chosen, and the projects contain archives from a wide range of domains with more than 4.3 million source software files. We divide the dataset into test and development sets in the ratio 70-20 on the directory level, and the

creation set is then randomly divided into validation and training sets in the ratio 85-15. The final executed algorithm is rebuilt using the complete dataset.

### 4.2.Preprocessing:

They handle the original code data in the following as a series of characters that correspond to a semantic analyzer's output. Interestingly, an in-depth traversal of a literal syntactic tree's final nodes can likewise be used to generate this (Central Standard Time). We avoid using the highest-level representations of structure in this study, such as concrete or conceptual syntax trees or flow control graphs, since they add extra complexity and constraints that hinder inference and decrease the breadth of the code-completing system. [13]Additionally, in the majority of programming languages, these models can only be accurately retrieved on syntactically correct entire code snippets something that code generation systems frequently do not have access to.

Several standardization criteria are taken from the program's concrete grammar tree and used for empirical language modeling of the source code, which is the foundation of our approach. Different styles and tab or white space traditions are an issue, so we use customized tokenizes to convert the source code into visual software tokens and then renew the text with a standard design. We parse each file's program code, gather data on term kinds, and use that data to organize the code, identify language related to sub-tokens, and compress the sequences throughout pre-processing. Inference and instruction are the two purposes for this.
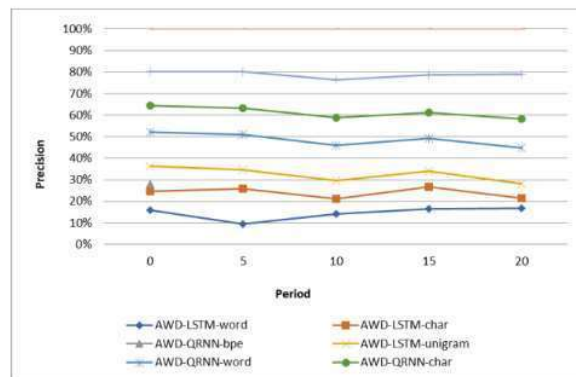
### 4.3.Training Results:

The outcomes obtained following the complete training of the chosen DNN architectures using the various tokenization methods are shown in this section. As mentioned in the preceding section, we used three different base approaches Word, Monogram, the use of BPE, and Char to train the AWD-LSTM and AWD-QRNN the deep neural network architecture utilizing distinct tokenization models.
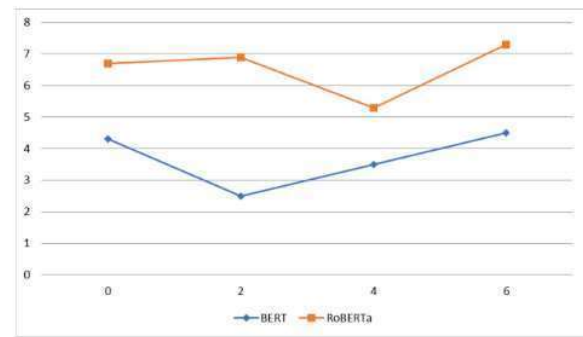
**Table 4.2:** The Ultimate Metrics for Each NN

| Deep Neural Network Architecture | Periods | Precision | Loss of Train | Delay in Validation | Already Trained? |
|---|---|---|---|---|---|
| The phrase AWD-LSTM | 29 | 0.493546 | 1.432564 | 2.984689 | Yes |
| AWD-LSTM monogram | 29 | 0.426576 | 2.546263 | 1.846438 | Yes |
| The AWD-LSTM BPE | 29 | 0.567325 | 2.426746 | 0.352155 | Yes |
| AWD-LSTM character | 29 | 0.425566 | 0.353516 | 1.423452 | Yes |
| Word with AWD-QRNN | 29 | 0.246742 | 0.345267 | 0.562664 | Yes |
| AWD-QRNN monogram | 29 | 0.746334 | 1.655265 | 2.512546 | Yes |
| The AWD-QRNN BPE | 29 | 0.976546 | 2.456356 | 1.678462 | Yes |
| AWD-QRNN character | 29 | 0.758493 | 2.256685 | 1.456577 | Yes |

For each AWD-LSTM and AWD-QRNN, we utilized one vintage of training to fit the simulation's head and thirty rounds of fine-tuning. To suit the head and refine the representations throughout ten periods, Transformer neural networks were trained evenly over one epoch.

**Figure 4.1:** Neural Networks' Precision in Producing Source Code



**Figure 4.2:** Development of neural network accuracy

They employed a two-pronged approach to assess the trained NNs: human output assessment of the algorithms and the application of NN learning measures. Reliability for the set of validations and cost for both the validation and training sets are two of the most widely used measures in research. By employing training and validation datasets, they aid researchers in comprehending the behavior of the neural network (NN), how the prediction is adapted to the information set, and its effectiveness and error scores. In this instance, reliability refers to the score indicating how well the LM can, given a string of words given the set of validations; anticipate the next word to fill in the missing ones. [14] Upon deploying the DNN to either the learning or verification set, the error is reported by the loss metrics. The GitHub source has all of the setup details for the deep neural networks and statistics. They evaluated the models' computational quality in addition to those criteria by using them in the suggested tasks to auto-complete and produce text and watch how well they worked. After training, Table 4.2 shows the ultimate metrics for each of the NNs. In a similar vein, Figures 4.1 and 4.2 illustrate how each model's precision changed after training.

On the one hand, the AWD-LSTM with charcoal tokenization when possible was the in general NN-tokenization simulation conjunction that carried out better in the case of precision metrics, based on the results shown in Table 4.1 and Figure 4.1. For neural network models intended for computerized source code generation, AWD-qRNN, and Transformer GPT-2. Conversely, the outcomes displayed in Table 4.1 and Figure 4.2 demonstrate that the two approaches (BERT and Roberta) achieved good accuracy results, with respective values of 0.985344 and 0.995632 in the transformer algorithms meant for predictive completion.

## 5.    CONCLUSION:

This work's suggested system can automatically produce compiled descriptions for C programs. Software code methods are the only ones for which description can be produced by current software description tools. NLP's code summary method is combined with natural language generation (NLG) in the autonomous software description system. The C program's entire text can be annotated by this technique. The application life cycle algorithm's documentation process is less expensive due to this automated source code annotation solution. The entire cost of developing a piece of software may also be decreased by automated processes. The outcomes from modeling languages used for creating code samples or auto-complete software features are compared concerning various methods to tokenization when possible designs, neural network designs, models that have been trained, and transfer learning. To find out which DNN architectures such as AWD LSTM, AWD-QRNN, and Transformer perform best with various encoding models word, monogram, BPE, and char we looked at a variety of deep neural network architectures. The precision was slightly lower than the

alternative architectures in larger models, such as Transformer GPT-2. On the original code creation tests, GPT-2 performed better, nevertheless.

## References:

[1] Arthur, M. P. (2020). Automatic source code documentation using code summarization technique of NLP. Procedia Computer Science, 171, 2522-2531.

[2] Le, T. H., Chen, H., & Babar, M. A. (2020). Deep learning for source code modeling and generation: Models, applications, and challenges. ACM Computing Surveys (CSUR), 53(3), 1-38.

[3] Ilame, N. (2024). Machine Learning-Powered Programming: Exploring the Fusion of AI and Coding. Innovative Computer Sciences Journal, 10(1), 1-11.

[4] Zhu, J., & Shen, M. (2020, April). Research on Deep learning Based Code generation from natural language Description. In 2020 IEEE 5th International Conference on Cloud Computing and Big Data Analytics (ICCCBDA) (pp. 188-193). IEEE.

[5] Maddigan, P., & Susnjak, T. (2023). Chat2vis: Generating data visualisations via natural language using chatgpt, codex and gpt-3 large language models. Ieee Access.

[6] Song, X., Sun, H., Wang, X., & Yan, J. (2019). A survey of automatic generation of source code comments: Algorithms and techniques. IEEE Access, 7, 111411-111428.

[7] Shim, S., Patil, P., Yadav, R. R., Shinde, A., & Devale, V. (2020, January). DeeperCoder: Code generation using machine learning. In 2020 10th Annual Computing and Communication Workshop and Conference (CCWC) (pp. 0194-0199). IEEE.

[8] Shin, J., & Nam, J. (2021). A survey of automatic code generation from natural language. Journal of Information Processing Systems, 17(3), 537-555.

[9] LeClair, A., Jiang, S., & McMillan, C. (2019, May). A neural model for generating natural language summaries of program subroutines. In 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE) (pp. 795-806). IEEE.

[10] Khurana, D., Koli, A., Khatter, K., & Singh, S. (2023). Natural language processing: State of the art, current trends and challenges. Multimedia tools and applications, 82(3), 3713-3744.

[11] Soliman, A., Shaheen, S., & Hadhoud, M. (2024). Leveraging pre-trained language models for code generation. Complex & Intelligent Systems, 1-26.

[12] Ye, W., Xie, R., Zhang, J., Hu, T., Wang, X., & Zhang, S. (2020, April). Leveraging code generation to improve code retrieval and summarization via dual learning. In Proceedings of The Web Conference 2020 (pp. 2309-2319).

[13] Svyatkovskiy, A., Deng, S. K., Fu, S., & Sundaresan, N. (2020, November). Intellicode compose: Code generation using transformer. In Proceedings of the 28th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering (pp. 1433-1443).

[14] Cruz-Benito, J., Vishwakarma, S., Martin-Fernandez, F., & Faro, I. (2021). Automated source code generation and auto-completion using deep learning: Comparing and discussing current language model-related approaches. AI, 2(1), 1-16.