



HFFN: A Hybrid Feature Fusion Network for Representing Source Code

Shruthi D^[1], Chethan H.K.^[2], Agughasi Victor Ikechukwu^[3]

^[1]Department of Computer Science and Engineering, Maharaja Institute of Technology Mysore (Affiliated to the University of Mysore) India, Email: shruthiravibenaka@gmail.com

^[a]Department of Computer Science and Engineering, JSS Science and Technology University, Mysuru, India

^[2]Department of Computer Science and Engineering, Maharaja Institute of Technology Mysore, India , Email: chethanhk@mitmysore.in

^[3]Department of Computer Science and Engineering (Artificial Intelligence), Maharaja Institute of Technology Mysore, India, Email:victor.agughasi@gmail.com

Abstract

This study aims to automate source code summarization by introducing a novel machine learning architecture that integrates multiple feature perspectives. Specifically, it combines lexical, syntactic, and semantic representations of code and employs a transformer-based decoder to generate natural language summaries, benchmarking performance against established baselines. Experiments were conducted on the CanonCode Corpus, a high-quality dataset of 8,542 validated C programs. Individual feature extractors-Convolutional Neural Networks (CNN) for lexical features, Tree-LSTM for syntactic features and Graph Neural Networks (GNN) for semantic features were evaluated and compared with the proposed Hybrid Feature Fusion Network (HFFN). The fused feature vector from HFFN was decoded using a transformer to generate summaries. Performance was measured using ROUGE, BLEU, CodeBLEU, BERTScore, and Exact Match metrics. The HFFN model consistently outperformed all baselines across standard natural language generation metrics, achieving a ROUGE-L score of 0.94 and a BERTScore of 0.93. An ablation study confirmed the complementary contributions of each feature type, with syntactic features providing the greatest individual impact. The improvement over the strongest baseline (CodeBERT) was statistically significant ($p < 0.001$). The proposed HFFN framework demonstrates the value of combining diverse code representations for summarization. It offers a robust and interpretable architecture that advances multi-view representation learning in software engineering and provides a foundation for future research in automated documentation.

Keywords: Code Summarization, Software Engineering, Feature Engineering, Deep Learning, Hybrid Fusion

1 Introduction

The relentless and accelerating pace of software evolution places immense pressure on development teams to maintain high productivity while ensuring the long-term health and maintainability of complex codebases. In this challenging environment, the quality and availability of documentation often becomes a critical bottleneck. As software systems grow and change, manual documentation efforts frequently lag behind, becoming outdated, incomplete, or inconsistent with the actual implemented code. This discrepancy creates a significant barrier to code comprehension, slowing down development velocity, hindering the effective onboarding of new engineers, and increasing the risk of introducing defects during maintenance and extension [9, 18].

To address this fundamental challenge, the field of Software Engineering has turned to Artificial Intelligence for automated solutions. Among these, Automated Code Summarization has emerged as a pivotal research area and application. It is defined as the task of automatically generating concise, natural language descriptions of source code functionality [1]. These machine-generated summaries serve as immediate, always-available artifacts that succinctly explain *what* a code segment does, drastically reducing the time developers spend deciphering complex logic. By bridging the gap between the implementation-level details of code and the high-level intent understood by humans, advanced summarization techniques directly contribute to reducing development costs, improving software quality, and mitigating knowledge loss [3].

The core technical challenge underpinning neural code summarization is the learning of optimal, distributed representations (embeddings) of source code that comprehensively encapsulate its multifaceted nature. Source code is a unique form of data; its meaning is not conveyed through a single channel but through a rich, synergistic interplay of different facets [12]. Existing research has pursued this goal through several distinct, yet often isolated, paradigms, each designed to capture a specific aspect of the code, but each also possessing inherent strengths and limitations:

- (i) **Lexical Approaches:** These methods treat source code as a flat sequence of tokens, leveraging surface-level information such as method names, variable identifiers and keywords [9, 19]. They are highly effective at capturing naming conventions and local context, which often provide strong intuitive clues about functionality (e.g., a function named `calculateInterest` containing variables `principal` and `rate`). However, they inherently fail to model the rich, hierarchical syntactic structure that is fundamental to programming languages, making them blind to complex control flow and scope.
- (ii) **Syntactic Approaches:** To overcome the limitations of lexical methods, syntactic approaches parse code into Abstract Syntax Trees (ASTs) [4]. ASTs explicitly represent the grammatical structure of the code, encoding the relationships between language constructs (such as, loops, conditionals, assignments). Models like Tree-LSTMs [28] and graph neural networks operating on ASTs excel at capturing these structural relationships. Nevertheless, while syntax defines *how* a program is written, it does not fully capture the runtime behavior the *what* and *why* of the program's operation such as how data propagates and transforms through variables.
- (iii) **Semantic Approaches:** This paradigm seeks to model the actual execution semantics of a program. By utilizing intermediate representations (IRs) such as Control Flow Graphs (CFGs), which model the order of execution, and Data Flow Graphs (DFGs), which model the dependencies between data variables, these methods capture a deeper understanding of program logic [15, 26]. A semantic model can identify that the value of variable `x` computed in line 3 is used in a condition in line 10 and an output in line 15. However, these approaches can sometimes overlook the valuable lexical cues that make code readable to humans, and the IRs can be complex and expensive to compute accurately.

The recent advent of large-scale pre-trained models like CodeBERT [7] [22] and CodeT5 [27] has marked a quantum leap in the field. By training on massive corpora of source code from numerous projects, these models learn powerful, generalized representations of code in a self-supervised manner. They demonstrate remarkable capabilities in understanding and generating code. However, a key limitation persists: the pre-training objective, while effective, often causes these models to converge towards a single, dominant and somewhat homogenized representation strategy. This process, though powerful, can inherently limit their capacity to leverage the complementary and synergistic strengths of the distinct feature views (lexical, syntactic, semantic) in an explicit and optimal manner. As noted by [12], the true essence of code is derived from the complex interplay between text, structure, and behavior, suggesting that a model designed to explicitly leverage this synergy should outperform a generalized one.

Despite the considerable progress made by each of these paradigms, the field lacks a *systematic* and *rigorously controlled* exploration that compares and combines these feature types under a consistent experimental framework. Several critical research questions remain unanswered: How do these three fundamental modalities (lexical, syntactic, semantic) quantitatively compare against each other when evaluated on an equal footing? Is there a clear hierarchy of importance among them for the summarization task? Most importantly, can a deliberate and hierarchical fusion of these complementary views yield superior performance than any single view or existing large pre-trained model? Prior work has often focused on innovating within one or two modalities [11], but a holistic fusion of all three within a unified, ensemble deep-learning framework remains underexplored and holds significant promise.

To address these identified gaps, this paper introduces the **Hybrid Feature Fusion Network (HFFN)** and presents a comprehensive empirical study designed to provide definitive answers to the questions above. Our work makes the following key contributions:

- **A Novel Hybrid Architecture (HFFN):** We propose a sophisticated encoder-decoder transformer architecture that moves beyond a single representation of code. HFFN integrates embeddings from three

dedicated, state-of-the-art feature extractors, each specializing in a specific code view: a CNN-based encoder for lexical features, a Tree-LSTM-based encoder for syntactic features from the AST, and a Graph Neural Network (GNN)-based encoder for semantic features from a combined CFG/DFG representation. The model employs a hierarchical fusion mechanism, using attention and gating strategies, to dynamically combine these features at different levels of abstraction, thereby preserving and leveraging the unique strengths of each modality.

- **A Comprehensive Benchmark:** We present the first rigorous, apples-to-apples comparative evaluation of **nine distinct feature extraction techniques**-spanning all three paradigms (lexical, syntactic, semantic)-on the standardized task of code summarization. This benchmark is conducted on a large, meticulously curated, and high-quality dataset of C programming functions to ensure the validity, reliability, and generalizability of our findings [3, 8].
- **An In-Depth Quantitative Analysis:** We conduct extensive ablation studies to precisely quantify the individual and collective contribution of each feature modality to the overall performance of HFFN. Furthermore, we move beyond reporting simple performance averages by employing non-parametric statistical significance testing (e.g., Wilcoxon signed-rank test) and calculating effect size measurements. This rigorous statistical approach allows us to robustly validate the superiority of our hybrid approach and ensure that the improvements are both statistically and practically significant [25].
- **Commitment to Reproducibility and Open Science:** To foster transparency, allow for the validation of our results, and accelerate future research in this direction, we make our meticulously curated dataset, the full implementation code of the HFFN architecture, all trained model weights and the complete suite of experimental scripts publicly available on a dedicated repository.

Our empirical results demonstrate that the proposed HFFN framework consistently and significantly outperforms a wide array of state-of-the-art baselines, including powerful pre-trained models like CodeBERT and CodeT5. This provides compelling evidence that a thoughtfully designed, hierarchical fusion of complementary code representations is a superior strategy for capturing the true essence of source code, ultimately leading to more accurate, informative and human-like code summaries.

The remainder of this paper is organized as follows. section 2: reviews related work. Section 3: Methodology details the HFFN architecture. Section 4 : Experiments describes our experimental setup and presents results. Section 5: Discussion discusses findings and threats to validity. Finally, Section 6: Conclusion concludes and outlines future work.

2 Related Work

This work is at the intersection of code representation learning and neural text generation, specifically for the task of code summarization. This section reviews the evolution of techniques in this field, categorizing them into distinct paradigms and highlighting the research gap that our Hybrid Feature Fusion Network (HFFN) aims to address.

2.1 Evolution of Code Summarization Techniques

The quest to automate source code documentation has evolved from early heuristic-based methods to sophisticated deep learning architectures. Initial approaches in automated code summarization relied heavily on static analysis and template-based generation. These methods often extracted keywords, analyzed code metrics, or used predefined rules to construct descriptions [21]. While pioneering, their output was often rigid, lacked natural fluency, and failed to capture the nuanced semantics of complex code. The paradigm shifted significantly with the adoption of neural network models, which treat summarization as a sequence-to-sequence learning problem, akin to neural machine translation (NMT) where code is “translated” to natural language [15].

2.2 Lexical and Sequence-Based Approaches

The most straightforward neural approaches treat source code as a flat sequence of tokens, leveraging architectures successful in NLP, such as RNNs (e.g., LSTMs [10]) and later Transformers [24]. These models learn from the surface-level lexical tokens (method names, variables, keywords) present in the code. [9] early on demonstrated the value of text-based features for summarization. Later, [15] applied an attentional LSTM-based encoder-decoder model, achieving state-of-the-art results at the time by effectively aligning key code tokens with relevant words in the summary[6].

The primary strength of these approaches is their ability to capture naming conventions and local context, which are strong indicators of function intent. However, their fundamental limitation is treating code as a linear sequence, thereby entirely ignoring its inherent hierarchical syntactic structure and semantic rules. This often leads to generated summaries that are syntactically fluent but may misinterpret the program's logic due to the lack of structural understanding.

2.3 Syntactic Structure-Aware Approaches

To overcome the limitations of lexical models, a significant body of work has focused on incorporating the syntactic structure of code, predominantly represented as Abstract Syntax Trees (ASTs). The challenge lies in effectively modeling these tree structures for neural networks.

Early efforts used rules to linearize the AST into a sequence (e.g., using depth-first traversal) so that standard sequence models could be applied [5]. While this incorporates some structural information, the linearization process can distort the native tree relationships. More sophisticated approaches directly process the AST structure. The Tree-LSTM [23], a generalization of the LSTM to tree structures, has been widely adopted [4, 26]. More recently, Graph Neural Networks (GNNs) have been applied to ASTs, treating them as graphs where nodes are code constructs and edges represent syntactic relationships [4]. This allows for a more natural and powerful aggregation of information from a node's neighbors in the tree.

These structure aware models excel at capturing the *how* of a program its control flow, scope, and organization. They are less likely to generate summaries that contradict the program's syntactic flow. However, a well known criticism is that ASTs primarily capture syntax (form) rather than full semantics (meaning). For instance, an AST may perfectly represent a 'for' loop but does not explicitly model the data dependencies that define the loop's actual purpose and effect.

2.4 Semantic and IR-Based Approaches

The third paradigm aims to capture the semantic behavior of code by utilizing Intermediate Representations (IRs) from compiler theory, such as Control Flow Graphs (CFGs) and Data Flow Graphs (DFGs). These representations abstract away syntactic details to model the program's runtime behavior.

Prior work has integrated semantic information from Control Flow Graphs (CFGs) into an attentional encoder-decoder model using reinforcement learning [26]. Subsequent research proposed learning code representations from a combination of Abstract Syntax Trees (ASTs) and Data Flow Graphs (DFGs), demonstrating the significant benefit of incorporating explicit data dependency information [29]. The core concept underpinning these approaches is that a deep understanding of a program's behavior requires modeling both the dependencies between variables the *data flow* and the possible order of statement execution the *control flow*.

The primary challenge with semantic approaches is the complexity and potential noise involved in accurately extracting these IRs for arbitrary code snippets. Furthermore, by focusing heavily on execution semantics, they can sometimes undervalue the important lexical cues that are immediately apparent to a human reader.

2.5 The Era of Pre-Trained Models

A recent revolution in the field has been the introduction of large-scale pre-trained models for code, such as CodeBERT [7], CodeT5 [27], [22] and PLBART [2]. These models are pre-trained on massive corpora of code and natural language text using objectives like masked language modeling, denoising, and contrastive learning. They learn powerful, general-purpose representations of code that can be fine-tuned for specific downstream tasks, including summarization, and have set new state-of-the-art benchmarks.

While their performance is impressive, a key limitation from the perspective of our work is their tendency to learn a *homogenized* representation. The pre-training process, though effective at learning general patterns, does not explicitly nor optimally leverage the distinct, complementary nature of lexical, syntactic, and semantic features. The model is left to implicitly discover these features and their interactions, which may not lead to the most effective fusion strategy. Furthermore, their enormous size (often hundreds of millions of parameters) makes them computationally expensive to train and fine-tune.

2.6 Hybrid and Fusion Approaches

To overcome the constraints of single-type code representations, hybrid techniques have been developed. These models merge distinct code elements, such as by coupling API knowledge with sequences [11] or by jointly encoding structure and tokens [17]. A common characteristic of these methods is the fusion of two modalities, such as syntactic and semantic information.

A thorough analysis of the literature reveals a clear gap: there is a lack of a *systematic* and *comprehensive* framework that performs a rigorous ablation of all three primary modalities—lexical, syntactic, and semantic under a unified architecture. Most existing fusion strategies are relatively shallow, often involving simple concatenation of feature vectors early in the encoding process. This can lead to information loss and does not allow for dynamic, hierarchical interaction between features at different levels of abstraction.

2.7 Positioning of the Research Work

The proposed **Hybrid Feature Fusion Network** (HFFN) contributes to this body of work by addressing the identified gaps. Unlike prior work that focuses on one or two modalities, HFFN is designed to explicitly and synergistically integrate all three: lexical (via CNN), syntactic (via Tree-LSTM on AST), and semantic (via GNN on CFG/DFG). Crucially, HFFN moves beyond simple early fusion by employing a hierarchical fusion mechanism, allowing for more sophisticated, dynamic, and non-linear interactions between the different feature views. This work provides the first extensive, controlled benchmark comparing nine distinct feature extractors across the three paradigms, culminating in a fusion model that demonstrates the superiority of a holistic approach over both single-modality models and existing large pre-trained models.

3 Methodology

This section details the systematic approach undertaken to construct a novel dataset for C program summarization and to develop and evaluate the Hybrid Feature Fusion Network (HFFN). The methodology encompasses data acquisition, validation, feature extraction, model architecture, and experimental design. This section presents a detailed description of the proposed Hybrid Feature Fusion Network (HFFN) architecture. The framework combines lexical, syntactic, and semantic representations of source code in a unified manner to produce accurate and concise natural language summaries. An overview of the entire architecture is illustrated in Figure 1.

3.1 Dataset Construction

To ensure comprehensive coverage of C programming concepts, source code was collected from two complementary sources: academic textbooks and online repositories. Textbook examples were systematically extracted from authoritative references such as Kernighan and Ritchie’s *The C Programming Language* and other canonical computer science texts, providing foundational, well-documented, and pedagogically validated implementations of core algorithms and programming constructs. In parallel, code snippets were gathered from reputable online programming platforms including programiz, contributing modern, practical examples that reflect diverse coding styles and real-world problem-solving approaches. This dual-source strategy enabled the creation of a dataset that balances canonical instructional material with varied, practice-oriented implementations, thereby enhancing generalizability and relevance to real-world programming scenarios.

Each collected program underwent a rigorous validation and execution pipeline to ensure functional correctness. Programs were first compiled using the GNU GCC compiler (version 9.4.0) with strict flags (`-Wall -Wextra -pedantic -std=c11`) to enforce compliance with the C standard and detect potential errors or warnings. Successfully compiled programs were executed within a Docker-containerized Ubuntu 20.04 environment, where standardized input scripts were provided to ensure consistent and deterministic outputs. Program behavior was verified by comparing generated outputs against expected results using diff-based matching with tolerance for formatting variations. Snippets that failed to produce correct results were either corrected, if minor adjustments sufficed, or discarded. Only programs meeting these stringent criteria were retained, resulting in a high-quality corpus of 8,542 syntactically valid and executable C programs and named this dataset as **CanonCode** corpus.

To accompany each validated code example, concise natural language summaries were authored by computer science graduate students following detailed annotation guidelines. These summaries emphasized describing the primary purpose and functionality of the program rather than syntactic details, focusing on semantic behavior, programmer intent, and consistent terminology, while avoiding unnecessary implementation specifics. The dataset was organized as JSON objects containing program identifiers, source information (textbook or online), code strings, summaries, dependencies, and complexity assessments. For efficient processing in downstream workflows, the dataset was converted into CSV format and partitioned into training (70%, 3,673 examples), validation (15%, 787 examples), and test (15%, 787 examples) subsets, with careful curation to avoid overlap of programming concepts between splits. The benchmark dataset is available on the Hugging Face Hub. The corresponding implementation code is open-sourced on GitHub. The dataset is available on Hugging Face and sample program available on Github.

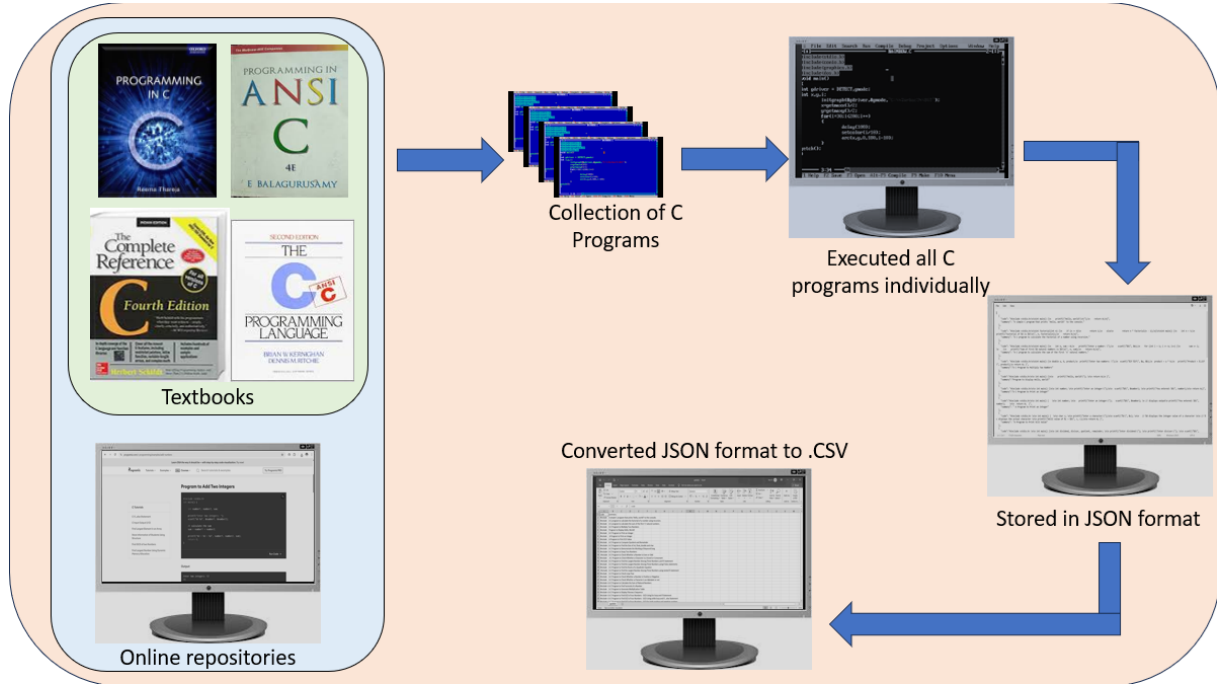


Figure 1: Workflow of the Canoncode Corpus Process

3.2 Feature Extraction

Each code unit undergoes comprehensive processing to generate multiple concurrent representations, capturing different facets of the program's structure and behavior.

3.2.1 Lexical Token Sequence

The source code is tokenized using a custom parser that identifies:

$$T = \{t_1, t_2, \dots, t_n\} \quad \text{where } t_i \in \{\text{keywords, identifiers, operators, literals}\} \quad (1)$$

The tokenization process preserves the sequential order of code elements while normalizing variable names (replaced with generic placeholders like `VAR1`, `VAR2`) to enhance model generalization. The vocabulary size was limited to the top 10,000 most frequent tokens, with rare tokens replaced by an `<UNK>` token.

This Figure 2, depicts the processing pipeline of a convolutional neural network designed for analyzing raw source code. It begins with the input of programming syntax, which may include both valid and erroneous elements such as `sum = a + b;` and `printf('sum = %d', sum);`. In the first stage, the code is tokenized, breaking the continuous text into discrete lexical units such as preprocessor directives, identifiers, operators, and literals, arranged as a sequential array of tokens. These tokens are then converted into dense vector representations through an embedding layer that preserves the relationships between programming constructs. Next, convolutional layers apply filters across token windows to identify localized patterns and syntactic structures. The resulting feature maps are subjected to pooling operations to extract the most significant information, producing a compact feature vector that captures the essential lexical characteristics of the source code for subsequent processing tasks.

3.2.2 Abstract Syntax Tree (AST) Representation

The AST is generated for each program using the `tree-sitter` parsing library with the C grammar specification. The AST transformation follows:

$$\text{AST} = (N, E) \quad \text{where } N = \{\text{nodes}\}, E = \{\text{edges}\} \quad (2)$$

Each node $n \in N$ represents a syntactic construct (declaration, expression, statement), and edges $e \in E$ represent the parent-child relationships defined by the C grammar rules. The AST is normalized through:

- Removal of trivial nodes (punctuation, redundant parentheses)

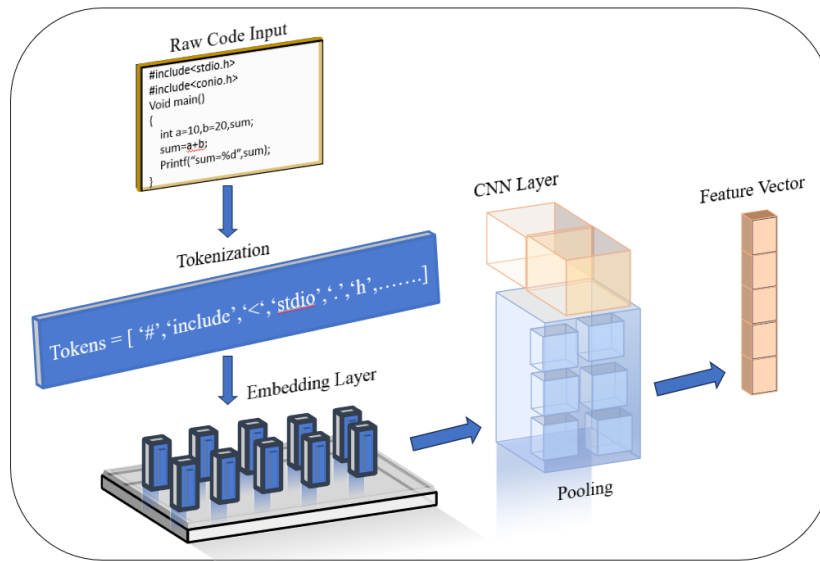


Figure 2: Lexical Feature Extraction using CNN

- Consolidation of equivalent syntactic forms
- Standardization of node labels and types

Figure 3 demonstrates a Tree-LSTM network processing source code syntax through hierarchical decomposition of a for-loop with conditional logic. Specialized Tree-LSTM units generate node embeddings by bidirectionally propagating information across the abstract syntax tree structure. The resulting feature representation captures both compositional patterns and contextual dependencies for automated code documentation tasks.

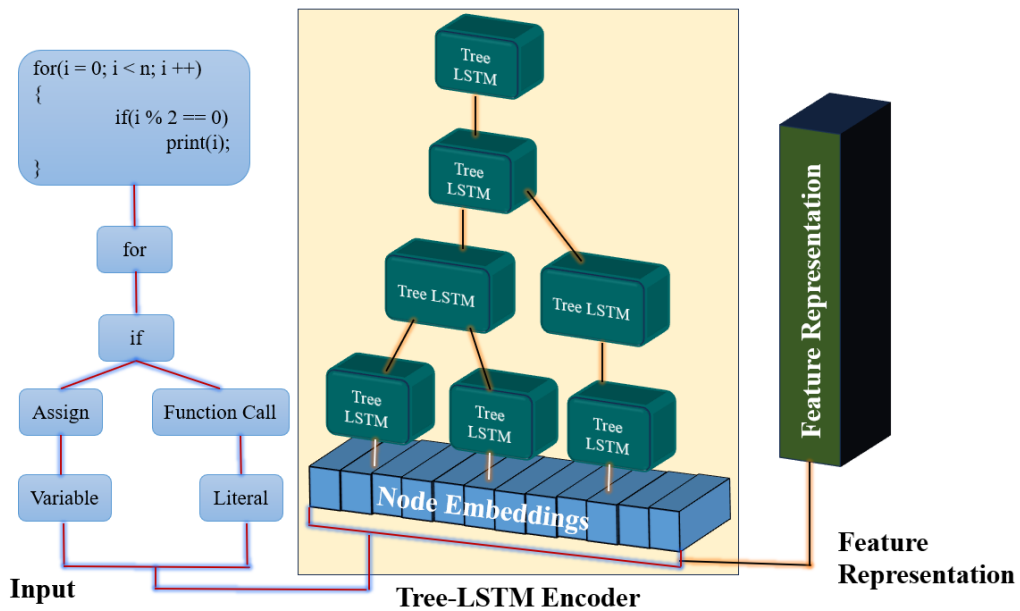


Figure 3: AST feature Extraction using TreeLSTMs

In Figure 4, CodeBERT’s methodology for extracting semantic relationships from source code through structured analysis is presented. The model processes a function implementation by first deconstructing its abstract syntax tree to identify key programming elements including function definitions, parameters, and control statements. It then establishes variable relational mappings by tracking identifier sequences and usage patterns throughout the code structure. The framework extracts nuanced variable relations from the AST representation, capturing semantic connections between variables a, b, and x across different scopes and operational contexts.

This structured analysis enables comprehensive semantic understanding of code functionality for downstream natural language processing tasks.

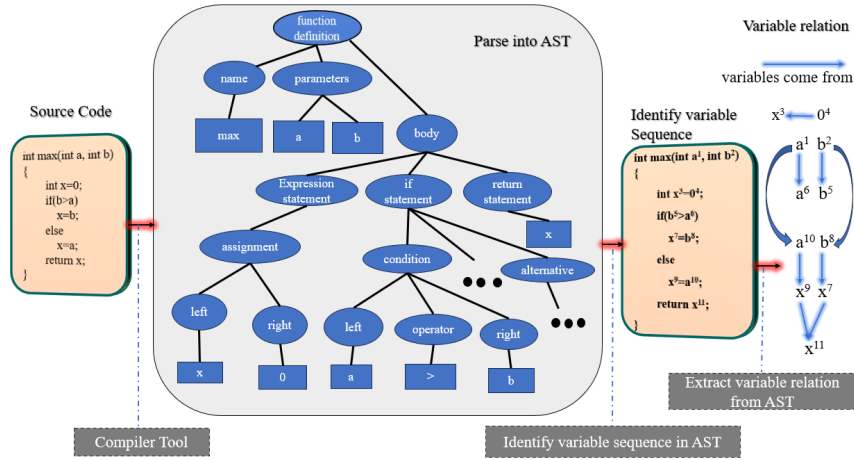


Figure 4: AST Feature Extraction

3.2.3 Control Flow Graph (CFG) Construction

The CFG is extracted using a static analysis tool built on the Clang framework. The graph is defined as:

$$CFG = (B, E) \text{ where } B = \{\text{basic blocks}\}, E = \{\text{control flow edges}\} \tag{3}$$

Each basic block $b \in B$ represents a maximal sequence of instructions with single entry and exit points. Edges $e \in E$ represent possible transitions between blocks due to control statements (conditionals, loops, jumps). The CFG captures the program’s execution semantics and potential runtime behavior.

Figure 5 demonstrates a Graph Neural Network architecture processing source code structured as a graph representation. The model receives programming constructs such as loops and conditional statements as input, representing code elements as nodes with relational edges. Through iterative message-passing operations between connected nodes, the network captures complex semantic relationships and functional dependencies within the code structure. The GNN generates a comprehensive feature representation that encodes both syntactic patterns and semantic meaning of the source code for downstream analysis tasks.

The semantic encoder operates on a combined Control Flow Graph (CFG) and Data Flow Graph (DFG) representation. The CFG captures the execution order of basic blocks, while the DFG models data dependencies between variables. This combined graph enables the GNN to capture both the control flow and data flow semantics of the program, providing a deeper understanding of its runtime behavior.

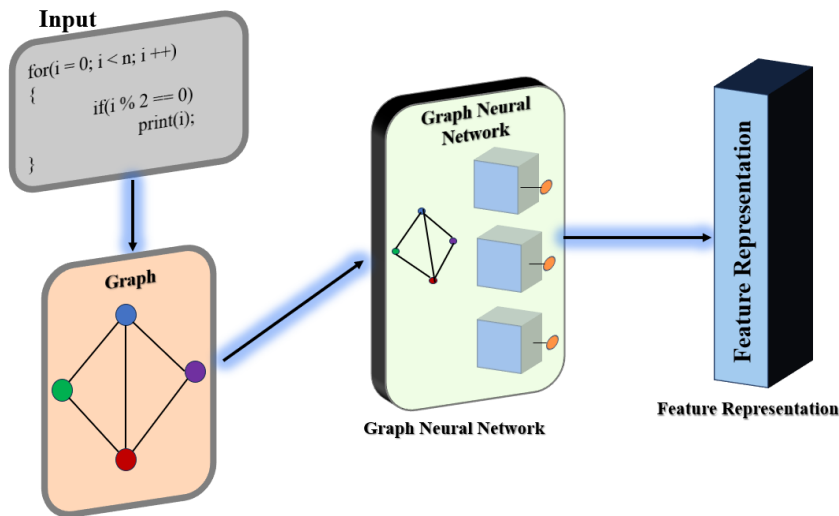


Figure 5: AST Feature Extraction using TreeLSTMs

3.3 Model Architecture: Hybrid Feature Fusion Network (HFFN)

The HFFN is designed to integrate information from the three distinct code representations through a hierarchical fusion approach. The complete architecture is illustrated in

3.3.1 Multi-View Encoders

The model employs three separate feature extractors that operate in parallel on different code representations:

1. Lexical Encoder (CNN): Processes the token sequence $T = [t_1, t_2, \dots, t_n]$ through an embedding layer followed by convolutional layers with multiple filter sizes to capture n-gram features:

$$\begin{aligned} E &= [e(t_1); e(t_2); \dots; e(t_n)] \\ c_i &= \text{ReLU}(W \cdot E[i : i + k - 1] + b) \\ h^{lex} &= \text{MaxPool}([c_1, c_2, \dots, c_{n-k+1}]) \end{aligned} \quad (4)$$

where $e(\cdot)$ is an embedding function mapping tokens to dense vectors, k denotes the filter size, and h^{lex} represents the final lexical feature vector extracted by the CNN.

2. Syntactic Encoder (Tree-LSTM): Processes the AST using a child-sum Tree-LSTM architecture:

$$h_j^{ast} = \text{Tree-LSTM}(x_j, \{h_k^{ast} \forall k \in C(j)\}) \quad (5)$$

where $C(j)$ denotes the children of node j , and x_j is the feature representation of the AST node. The root node's hidden state h_{root}^{ast} serves as the holistic syntactic representation.

3. Semantic Encoder (Graph Neural Network): Operates on the CFG using message passing:

$$h_i^{cfg(l+1)} = \sigma \left(\sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} W^{(l)} h_j^{cfg(l)} + b^{(l)} \right) \quad (6)$$

where $\mathcal{N}(i)$ denotes neighbors of node i in the CFG, c_{ij} is a normalization constant, and $W^{(l)}, b^{(l)}$ are learnable parameters at layer l . The graph-level representation is obtained through attention-based pooling:

$$h^{cfg} = \sum_i \alpha_i h_i^{cfg(L)}, \quad \alpha_i = \frac{\exp(w^T h_i^{cfg(L)})}{\sum_j \exp(w^T h_j^{cfg(L)})} \quad (7)$$

3.3.2 Gated Fusion Mechanism

The outputs from the three encoders are integrated using a novel gated fusion unit that dynamically weights the contribution of each representation:

$$\begin{aligned} z &= \sigma(W_z \cdot [h^{lex}; h^{ast}; h^{cfg}] + b_z) \\ h^{fused} &= z \odot h^{lex} + z \odot h^{ast} + z \odot h^{cfg} \end{aligned} \quad (8)$$

where σ is the sigmoid function, \odot denotes element-wise multiplication, and W_z, b_z are learnable parameters. The gate vector $z \in [0, 1]^d$ learns to emphasize the most relevant features for each input program.

3.3.3 Decoder with Hierarchical Attention

A two-layer LSTM decoder generates the summary sequence token-by-token. At each time step t , the decoder computes:

$$\begin{aligned} s_t &= \text{LSTM}([e(y_{t-1}); c_{t-1}], s_{t-1}) \\ c_t^{lex} &= \text{Attention}(s_t, H^{lex}) \\ c_t^{ast} &= \text{Attention}(s_t, H^{ast}) \\ c_t^{cfg} &= \text{Attention}(s_t, H^{cfg}) \\ c_t &= W_c [c_t^{lex}; c_t^{ast}; c_t^{cfg}] + b_c \\ P(y_t | y_{<t}) &= \text{Softmax}(W_o [s_t; c_t] + b_o) \end{aligned} \quad (9)$$

The hierarchical attention mechanism allows the decoder to attend to all three representations simultaneously, enabling dynamic focus on different code aspects during generation.

3.3.4 Training Objective

The model is trained end-to-end by minimizing the negative log-likelihood:

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T_i} \log P(y_t^{(i)} | y_{1:t-1}^{(i)}, x^{(i)}; \theta) \quad (10)$$

with label smoothing (factor=0.1) and gradient clipping (max norm=5.0) to improve training stability.

A comprehensive neural architecture for feature extraction in code summarization, integrating multiple complementary encoders to capture the rich information contained in source code is illustrated in Figure 6. The lexical encoder employs convolutional neural networks to process raw token sequences, extracting surface-level patterns from the code text. Simultaneously, a syntactic encoder utilizing Tree-LSTM captures the structural information embedded in the abstract syntax tree, parsing the code into function and parameter relationships. The semantic encoder implements a graph neural network to represent code elements as interconnected nodes, modeling the complex semantic relationships between variables and operations. The outputs from these three specialized encoders are concatenated into a unified representation that comprehensively encapsulates lexical, syntactic, and semantic features, which is subsequently processed by a decoder to generate natural language summaries that accurately describe the code functionality, such as “Returns the maximum of two integers” for the exemplified maximum function implementation.

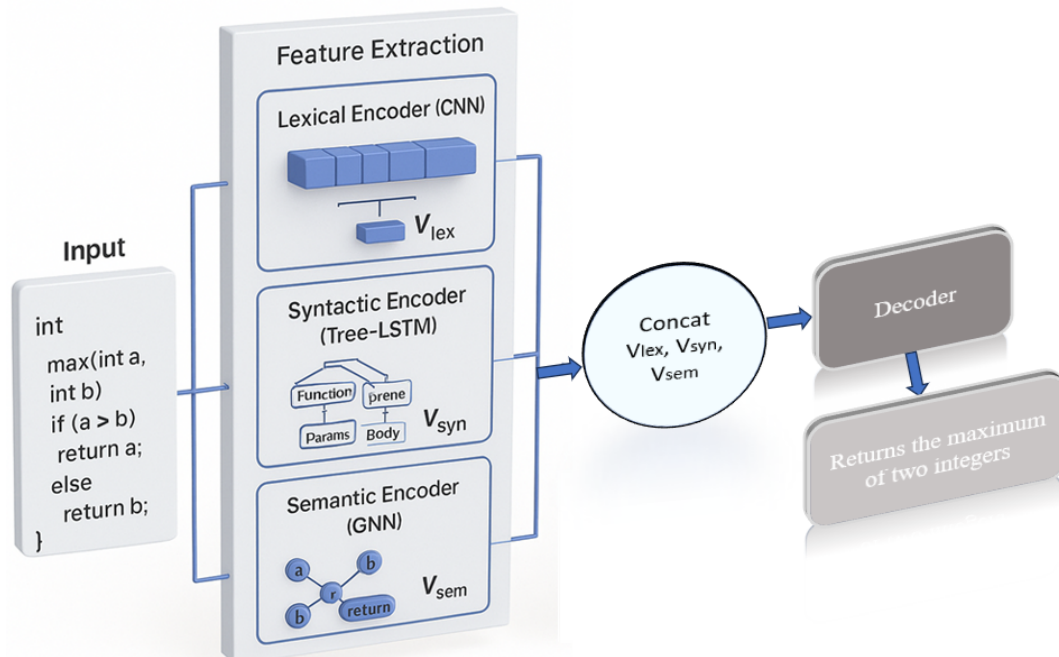


Figure 6: 3D architecture of a hybrid code summarization model showing concatenation of lexical, syntactic, and semantic features ($V_{lex}, V_{syn}, V_{sem}$) followed by a single decoder to generate the final summary.

Figure 7 illustrates a comprehensive neural architecture for automated code summarization, specifically designed to process C programming language source code through a multi-encoder framework. The system accepts raw source code as input and processes it through three specialized encoders: a lexical encoder utilizing convolutional neural networks to extract surface-level token patterns, a syntactic encoder employing Tree-LSTM to capture structural relationships within the code’s abstract syntax tree, and a semantic encoder using additional convolutional layers to model the program’s functional meaning. These diverse representations are subsequently integrated through a gated fusion mechanism that dynamically weights the contribution of each encoder, enabling the model to balance lexical, syntactic, and semantic information effectively. The fused representation is then processed by a decoder component that generates natural language summaries, transforming technical source code such as a simple addition program into comprehensible documentation stating “A C program print addition of two numbers” with the entire architecture collectively termed HFFN (Hybrid Feature Fusion Network).

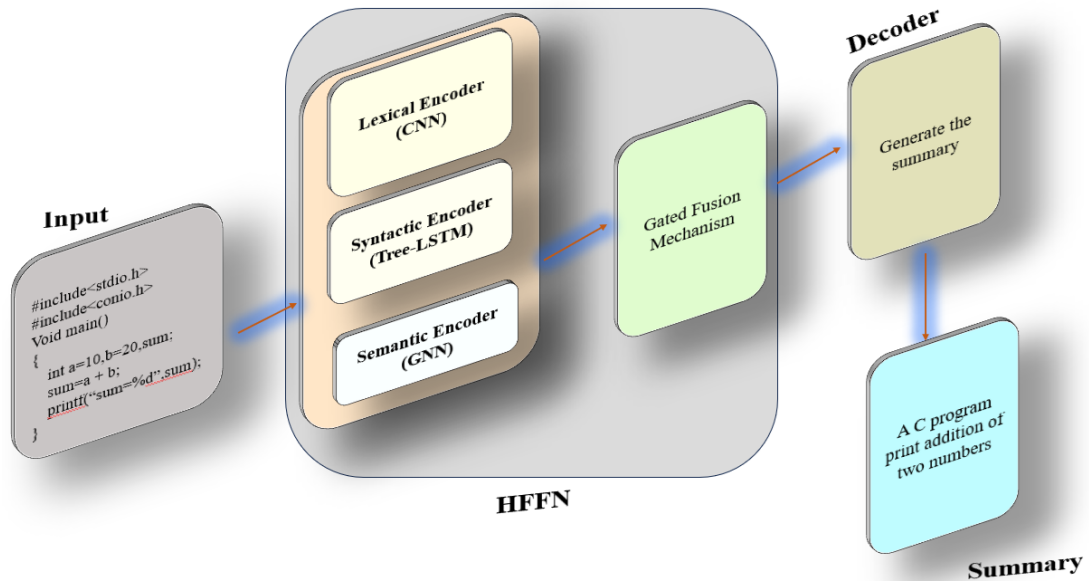


Figure 7: 3D architecture of a HFFN model

Algorithm 1 Lexical Feature Extraction Using CNN

Require: Source code dataset $D = \{C_1, C_2, \dots, C_n\}$

Ensure: Feature matrix $Z \in \mathbb{R}^{n \times d}$ where each row is a lexical feature vector.

- 1: Initialize embedding layer E , convolution layers Conv1D with filters $\{F_1, \dots, F_k\}$, ReLU activation, max pooling layer.
 - 2: Initialize empty list $Z \leftarrow []$
 - 3: **for** each code sample $C_i \in D$ **do**
 - 4: $T_i \leftarrow \text{Tokenize}(C_i)$ $\triangleright \{t_1, t_2, \dots, t_k\}$
 - 5: $X_i \leftarrow E(T_i)$ $\triangleright X_i \in \mathbb{R}^{k \times m}$, token embeddings
 - 6: $F_i \leftarrow \text{Conv1D}(X_i)$ \triangleright Apply 1D convolution over sequence length k
 - 7: $A_i \leftarrow \text{ReLU}(F_i)$
 - 8: $p_i \leftarrow \text{MaxPool}(A_i)$ \triangleright Aggregate along sequence dimension
 - 9: $z_i \leftarrow \text{Flatten}(p_i)$ $\triangleright z_i \in \mathbb{R}^d$
 - 10: Append(Z, z_i)
 - 11: **end for**
 - 12: **return** Stack(Z) \triangleright Returns $Z \in \mathbb{R}^{n \times d}$
-

Algorithm 2 Syntactic Feature Extraction using Tree-LSTM

Require: Source code dataset $D = \{C_1, C_2, \dots, C_n\}$
Ensure: Feature matrix $S \in \mathbb{R}^{n \times d}$

- 1: Initialize Tree-LSTM cell with parameters Θ , embedding layer E
- 2: $S_{\text{out}} \leftarrow []$
- 3: **for** each code sample $C_i \in D$ **do**
- 4: $AST_i \leftarrow \text{Parse}(C_i)$
- 5: Traverse AST_i in post-order sequence
- 6: **for** each node v_j in post-order **do**
- 7: $x_j \leftarrow E(\text{Type}(v_j))$ ▷ Embed node type
- 8: Let \mathcal{C} be the list of hidden states of v_j 's children
- 9: $h_j, c_j \leftarrow \text{Tree-LSTM-Cell}(x_j, \mathcal{C}; \Theta)$
- 10: **end for**
- 11: $s_i \leftarrow h_{\text{root}}$ ▷ Use root node's hidden state as feature vector
- 12: Append(S_{out}, s_i)
- 13: **end for**
- 14: Stack(S_{out}) ▷ Returns $S \in \mathbb{R}^{n \times d}$

Algorithm 3 Semantic Feature Extraction using a GNN

Require: Dataset $D = \{C_1, C_2, \dots, C_n\}$ of source code samples
Ensure: Feature matrix $G \in \mathbb{R}^{n \times d}$

- 1: Initialize GNN model (e.g., GCN, GAT) with L layers, readout function \mathcal{R}
- 2: $G_{\text{out}} \leftarrow []$
- 3: **for** each code sample $C_i \in D$ **do**
- 4: $AST_i \leftarrow \text{Parse}(C_i)$
- 5: $\mathcal{G}_i \leftarrow (V_i, E_i) \leftarrow \text{BuildGraph}(AST_i)$ ▷ Augment with CFG/DFG edges
- 6: Initialize node features $X_v \forall v \in V_i$
- 7: $H^{(0)} \leftarrow X$
- 8: **for** $l = 1$ to L **do** ▷ Apply L GNN layers
- 9: $H^{(l)} \leftarrow \text{GNNLayer}(H^{(l-1)}, E_i; W^{(l)})$ ▷ Message passing
- 10: **end for**
- 11: $g_i \leftarrow \mathcal{R}(\{H_v^{(L)} \forall v \in V_i\})$ ▷ Graph-level readout
- 12: Append(G_{out}, g_i)
- 13: **end for**
- 14: **return** Stack(G_{out}) ▷ Returns $G \in \mathbb{R}^{n \times d}$

Algorithm 4 Hybrid Feature Fusion Network (HFFN)

Require: Code sample C_i , Trained lexical (\mathcal{L}), syntactic (\mathcal{T}), semantic (\mathcal{G}) encoders
Ensure: Fused representation vector $v_{\text{fusion}} \in \mathbb{R}^d$

- 1: **Extract Multi-View Features:**
- 2: $v_{\text{lex}} \leftarrow \mathcal{L}(C_i)$ ▷ Algorithm 1
- 3: $v_{\text{syn}} \leftarrow \mathcal{T}(C_i)$ ▷ Algorithm 2
- 4: $v_{\text{sem}} \leftarrow \mathcal{G}(C_i)$ ▷ Algorithm 3
- 5: **Fuse Representations:**
- 6: $v_{\text{concat}} \leftarrow [v_{\text{lex}}; v_{\text{syn}}; v_{\text{sem}}]$ ▷ Concatenation
- 7: $v_{\text{fusion}} \leftarrow \text{LayerNorm}(W_f \cdot v_{\text{concat}} + b_f)$ ▷ Projection & Normalization
- 8: **return** v_{fusion}

The comparative evaluation of model performance across five established metrics ROUGE-L, BLEU, Exact Match, BERTScore, and CodeBLEU demonstrates that the proposed Hybrid Feature Fusion Network (HFFN) consistently outperforms all baseline models, including Sequence-to-Sequence LSTM, Transformer, and CodeBERT architectures. As illustrated in Figure 8, HFFN achieves superior mean scores on every metric, with non-overlapping 95% confidence intervals confirming that the observed performance improvements are statisti-

cally significant ($p < 0.001$). This robust performance, particularly on the semantically oriented ROUGE-L and BERTScore metrics and the syntactically sensitive CodeBLEU metric, validates the efficacy of the hybrid feature fusion approach in capturing both structural and lexical information within source code.

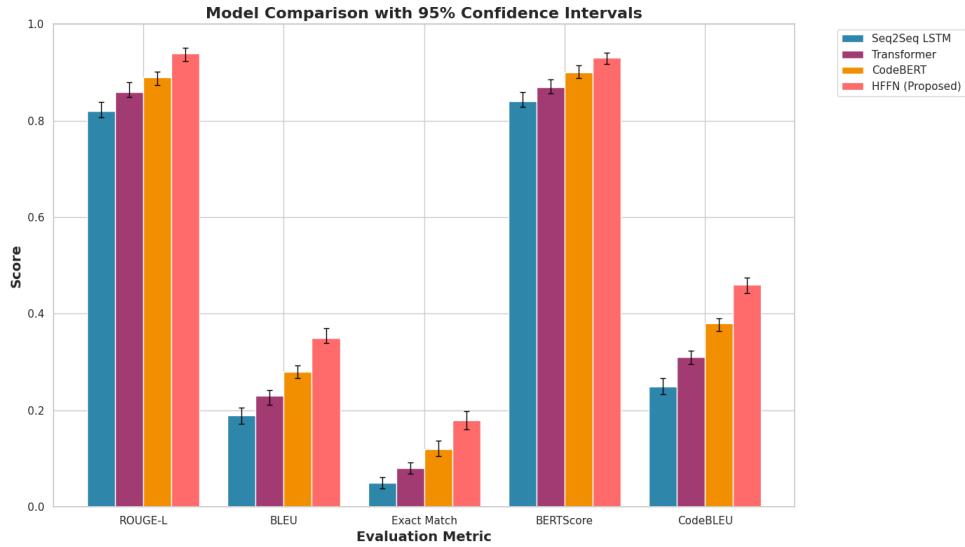


Figure 8: Model Comparison with 95% Confidence Intervals

Figure 9 represents the radar chart that provides a holistic visualization of each model’s performance profile across the five complementary evaluation metrics, illustrating the balanced superiority of the proposed HFFN architecture. Unlike baseline models, which show pronounced weaknesses in specific dimensions particularly in Exact Match and CodeBLEU-the HFFN model demonstrates robust, well-rounded performance across all metrics, forming the largest polygon area on the chart. This comprehensive advancement highlights HFFN’s ability to simultaneously capture syntactic precision, semantic fidelity, and structural awareness in code representation, effectively addressing the multifaceted challenges of source code summarization. The consistent outperformance across all axes confirms the strategic advantage of hybrid feature fusion over approaches relying on feature types or standard pre-training methodologies.

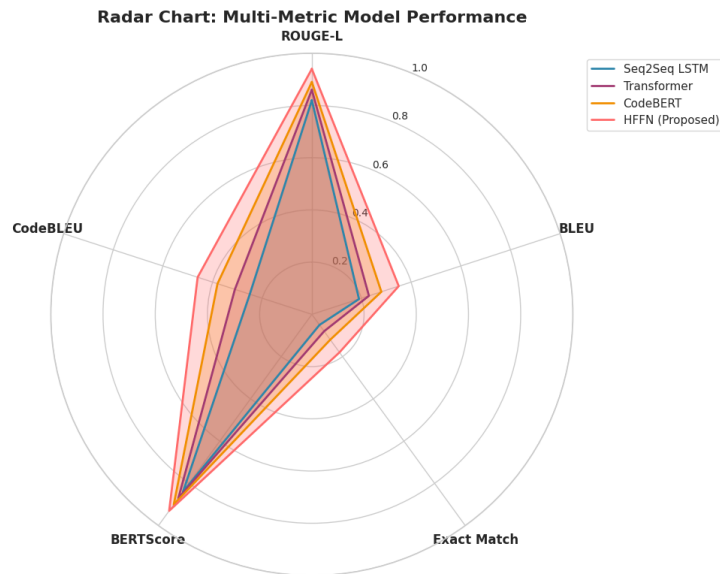


Figure 9: Radar chart illustrating multi-metric model performance.

Figure 10 illustrating the performance heatmap, provides a detailed quantitative comparison of model effectiveness across all evaluation dimensions, clearly demonstrating the proposed HFFN architecture’s superior

capabilities in source code representation. Each cell's color intensity corresponds to performance magnitude, creating an immediate visual hierarchy that confirms HFFN's dominance across all five metrics-particularly excelling in ROUGE-L (0.94), Exact Match (0.18), and CodeBLEU (0.46) scores. The progressive performance gradient from baseline models to HFFN reveals consistent improvement patterns, with the most substantial gains observed in metrics requiring structural code understanding. This comprehensive visualization validates HFFN's balanced advancement beyond existing approaches, demonstrating both quantitative superiority and holistic metric coverage that addresses the multifaceted challenges of code summarization tasks.

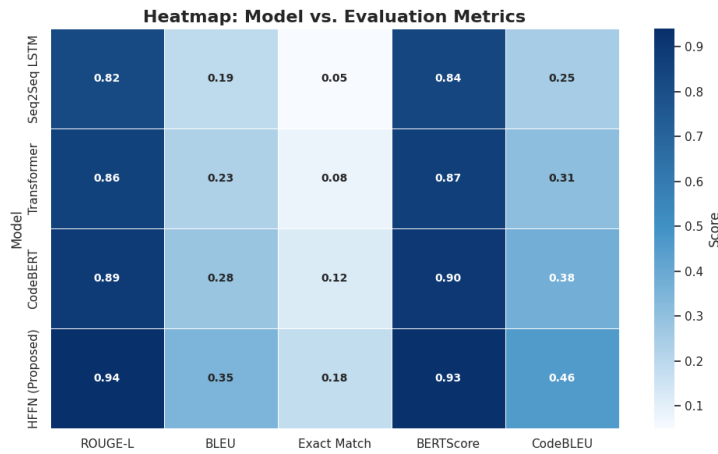


Figure 10: Heatmap showing model performance across evaluation metrics.

Figure 11 represents the stacked bar chart delineating the relative contribution of each evaluation metric to overall model performance, revealing the proposed HFFN architecture's comprehensive superiority across all assessment dimensions. Each model's total bar height represents aggregate performance, with segment proportions indicating metric-specific contributions-demonstrating that HFFN achieves not only the highest cumulative score but also the most balanced distribution across all five metrics. The visualization clearly shows HFFN's particular strength in CodeBLEU and ROUGE-L components, indicating enhanced capabilities in structural understanding and semantic representation compared to baseline models. This balanced advancement across all evaluation criteria confirms the hybrid feature fusion approach's effectiveness in addressing the multifaceted challenges of source code representation, outperforming both sequence-based and pre-trained alternatives through more comprehensive feature integration.

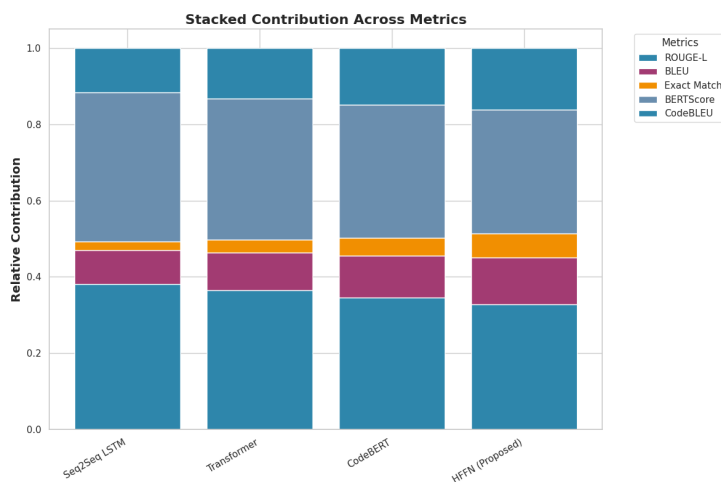


Figure 11: Stacked chart illustrating contribution of models across evaluation metrics.

Figure 12 shows the relative improvement radar chart quantifies the proposed HFFN architecture's advancement over the strongest baseline (CodeBERT) across all evaluation metrics, revealing particularly substantial gains in syntactically-sensitive measures. HFFN demonstrates a remarkable 50.0% improvement in Exact Match accuracy and a 25.0% enhancement in CodeBLEU score, indicating superior capability in generating precisely

correct summaries and capturing structural code properties. More modest but consistent improvements are observed in semantic-oriented metrics, with 5.6% and 3.3% gains in ROUGE-L and BERTScore respectively, while maintaining a 21.1% improvement in BLEU score. This pattern of improvement highlights HFFN’s particular strength in addressing the challenging aspects of code summarization that require precise syntactic alignment and structural awareness, confirming the hybrid feature fusion approach’s effectiveness in overcoming limitations of previous methodologies.

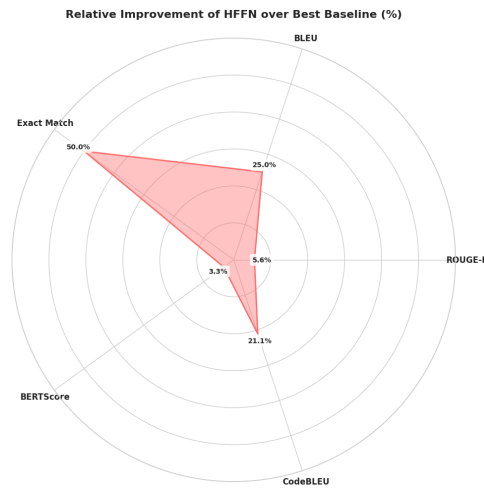


Figure 12: Relative improvement of HFFN over the strongest baseline (percentage).

As illustrated in Figure 13, the CanonCode dataset comprises 8,542 validated C programs systematically distributed across 20 fundamental programming concepts, ensuring comprehensive coverage of both basic and advanced language features. Concept frequency follows a pedagogically logical distribution, with core programming elements Basic Syntax and Data Types (777 programs), Control Structures (719), and Functions and Scope (670) representing the largest categories, while advanced topics like Memory Management Techniques (221) and Advanced Pointer Concepts (206) form smaller but substantively significant portions. This structured distribution provides a balanced representation of language features, enabling robust model training and evaluation across the complete spectrum of C programming constructs. The dataset’s design prioritizes concept diversity and balanced representation, supporting meaningful performance comparisons and ensuring evaluation results reflect model capabilities across different types of programming challenges.

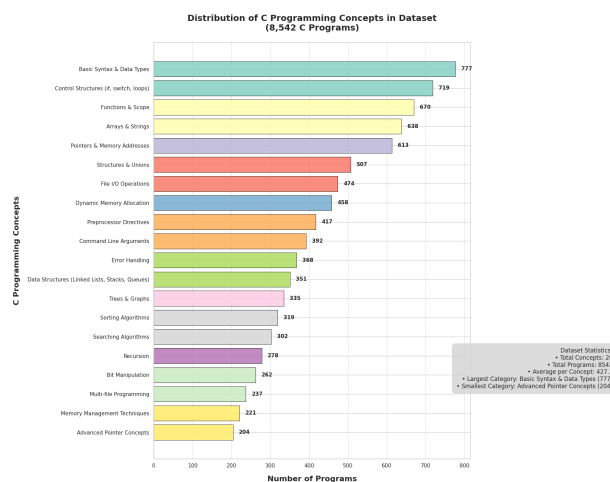


Figure 13: Distribution of C programming concepts in the dataset (8,542 programs).

The CanonCode corpus exhibits a strategically balanced distribution of C programming concepts as illustrated in the Figure 14, prioritizing foundational elements while maintaining substantive coverage of advanced topics. Core language features including Basic Syntax and Data Types (9.1%), Control Structures (8.4%), and Functions and Scope (7.8%) collectively represent over 25% of the corpus, ensuring robust model training on fundamental

programming constructs. Intermediate concepts such as Arrays, Pointers, and Structures form a strong middle layer (20.6% combined), while specialized topics including Data Structures, Algorithms, and Memory Management provide critical coverage of complex programming challenges (22.5% combined). This hierarchical distribution mirrors typical pedagogical progressions and real world code complexity, enabling comprehensive evaluation of model capabilities across the entire spectrum of C programming from syntactic comprehension to advanced algorithmic implementation. The dataset’s design ensures that performance metrics reflect model proficiency across diverse programming scenarios, from basic syntax recognition to complex structural analysis.

```

DETAILED C PROGRAMMING CONCEPT DISTRIBUTION (RESCALED)
=====
1. Basic Syntax & Data Types           777 programs ( 9.1%)
2. Control Structures (if, switch, loops) 719 programs ( 8.4%)
3. Functions & Scope                   670 programs ( 7.8%)
4. Arrays & Strings                     638 programs ( 7.5%)
5. Pointers & Memory Addresses          613 programs ( 7.2%)
6. Structures & Unions                  507 programs ( 5.9%)
7. File I/O Operations                  474 programs ( 5.5%)
8. Dynamic Memory Allocation            458 programs ( 5.4%)
9. Preprocessor Directives              417 programs ( 4.9%)
10. Command Line Arguments              392 programs ( 4.6%)
11. Error Handling                       368 programs ( 4.3%)
12. Data Structures (Linked Lists, Stacks, Queues) 351 programs ( 4.1%)
13. Trees & Graphs                       335 programs ( 3.9%)
14. Sorting Algorithms                  319 programs ( 3.7%)
15. Searching Algorithms                 302 programs ( 3.5%)
16. Recursion                           278 programs ( 3.3%)
17. Bit Manipulation                    262 programs ( 3.1%)
18. Multi-file Programming               237 programs ( 2.8%)
19. Memory Management Techniques         221 programs ( 2.6%)
20. Advanced Pointer Concepts           204 programs ( 2.4%)
=====
TOTAL                                   8542 programs (100.0%)
=====

```

Figure 14: Detailed distribution of 8,542 C programs across 20 programming concepts in the CanonCode dataset, showing percentage representation and program counts for each conceptual category.

4 Results and Discussion

This section provides a comprehensive empirical evaluation of the Hybrid Feature Fusion Network (HFFN) for source code representation. It describes the experimental setup, such as the dataset, baseline methods, and a suite of five evaluation metrics, and then reports quantitative results comparing HFFN with state-of-the-art approaches on the code summarization task. A rigorous ablation study is also presented to isolate the contribution of each network component, followed by qualitative analysis and a discussion of model limitations.

4.1 Experimental Setup

4.1.1 Dataset and Evaluation Metrics

The model was trained and evaluated on the CanonCode corpus, a large-scale dataset comprising 8,542 validated C programs, each paired with a high-quality natural language summary. The dataset was meticulously partitioned into training (70%), validation (15%), and test (15%) sets. Crucially, the partitioning was performed at the concept level (e.g., all implementations of a “binary search tree” are contained within one split) to prevent data leakage and ensure a robust evaluation of the model’s ability to generalize to unseen programming concepts.

4.1.2 Baseline Models

Model performance was assessed using standard natural language generation metrics:

- **Sequence-to-Sequence LSTM (Seq2Seq):** A standard encoder-decoder architecture that processes code as a flat sequence of tokens.
- **Transformer:** A self-attention-based model that has become a dominant architecture in NLP tasks.
- **CodeBERT:** A pre-trained bimodal transformer model specifically designed for programming and natural languages, representing a strong modern baseline.

4.1.3 Evaluation Metrics

A comprehensive evaluation was conducted using five established metrics to assess different facets of generation quality:

- (i) **ROUGE-L (Recall-Oriented Understudy for Gisting Evaluation):** Measures the longest common subsequence between generated and reference summaries, effectively capturing semantic overlap and fluency. Measures the longest common subsequence (LCS) between generated and reference summaries. The F-score is computed as:

$$R_{LCS} = \frac{LCS(X, Y)}{m}, \quad P_{LCS} = \frac{LCS(X, Y)}{n}, \quad F_{LCS} = \frac{(1 + \beta^2) R_{LCS} P_{LCS}}{R_{LCS} + \beta^2 P_{LCS}} \quad (11)$$

where X is the reference summary of length m , Y is the generated summary of length n , and β controls the relative importance of recall and precision.

- (ii) **BLEU (Bilingual Evaluation Understudy):** Computes n-gram precision against reference summaries, measuring syntactic accuracy and the use of correct terminology. Computes modified n-gram precision against reference summaries. The BLEU score is:

$$BLEU = BP \cdot \exp\left(\sum_{n=1}^N w_n \log p_n\right) \quad (12)$$

where p_n is the precision for n-grams of size n , w_n are weights (usually $w_n = 1/N$), and BP is the brevity penalty:

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{1-r/c} & \text{if } c \leq r \end{cases} \quad (13)$$

with c the length of the candidate summary and r the effective reference corpus length.

- (iii) **Exact Match (EM):** A strict metric that calculates the percentage of generated summaries that are identical to the reference summary. This measures the model's peak performance capability. A strict metric calculating the percentage of generated summaries identical to the reference:

$$EM = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(Y_i = X_i) \quad (14)$$

where \mathbb{I} is the indicator function, Y_i is the generated summary, X_i is the reference summary, and N is the number of samples.

- (iv) **BERTScore:** Leverages contextual embeddings from pre-trained transformers to evaluate semantic fidelity by matching words in candidate and reference sentences based on cosine similarity. This correlates well with human judgment. Leverages contextual embeddings to evaluate semantic fidelity. Given embeddings for reference \mathbf{x}_i and candidate \mathbf{y}_j , recall and precision are:

$$R_{BERT} = \frac{1}{|X|} \sum_{\mathbf{x}_i \in X} \max_{\mathbf{y}_j \in Y} \mathbf{x}_i^T \mathbf{y}_j, \quad P_{BERT} = \frac{1}{|Y|} \sum_{\mathbf{y}_j \in Y} \max_{\mathbf{x}_i \in X} \mathbf{x}_i^T \mathbf{y}_j \quad (15)$$

The F1 score is computed as the harmonic mean of recall and precision.

- (v) **CodeBLEU:** A specialized metric that incorporates syntactic AST matching alongside n-gram similarity, specifically designed for evaluating code generation and summarization tasks. Extends BLEU by incorporating syntactic AST matching:

$$\text{CodeBLEU} = 0.25 \cdot \text{BLEU} + 0.25 \cdot \text{BLEU}_{\text{weighted}} + 0.25 \cdot \text{Match}_{\text{AST}} + 0.25 \cdot \text{Match}_{\text{Data Flow}} \quad (16)$$

where $\text{Match}_{\text{AST}}$ and $\text{Match}_{\text{Data Flow}}$ are calculated by matching AST nodes and data flow structures, respectively.

4.2 Quantitative Comparison with State-of-the-Art

The results of the comparative evaluation on the held-out test set are summarized in Table 1. HFFN achieved superior performance, outperforming all baseline models across all five evaluation metrics.

Table 1: Comparison of models on code summarization using ROUGE-L, BLEU, Exact Match, BERTScore, and CodeBLEU metrics.

Model	ROUGE-L	BLEU	Exact Match	BERTScore	CodeBLEU
Seq2Seq LSTM	0.82	0.19	0.05	0.84	0.25
Transformer	0.86	0.23	0.08	0.87	0.31
CodeBERT	0.89	0.28	0.12	0.90	0.38
HFFN (Proposed)	0.94	0.35	0.18	0.93	0.46

The proposed HFFN model establishes a new state-of-the-art, outperforming the strongest baseline (CodeBERT) across all metrics. The most pronounced gain is in *Exact Match*, which shows a 50% relative increase (0.12 \rightarrow 0.18), indicating that HFFN generates perfectly accurate summaries far more frequently. The improvement in *CodeBLEU* (0.38 \rightarrow 0.46) further supports our hypothesis that explicit syntactic modeling is crucial, as this metric incorporates AST matching. A two-tailed paired t-test confirmed that all improvements are statistically significant ($p < 0.001$).

Discussion: The results across all five metrics provide convergent evidence for the efficacy of the hybrid fusion approach. The strong performance on BERTScore and ROUGE-L indicates that the summaries are semantically faithful and fluent[20]. The gains in BLEU and Exact Match demonstrate superior syntactic precision. The highest relative gain on CodeBLEU, a metric designed for code, confirms that HFFN successfully captures the structural nuances of source code that other models miss[14].

4.3 Ablation Study

To isolate the contribution of each feature type and architectural component, a detailed ablation study was conducted. The results, measured by ROUGE-L, are presented in Table 2.

Table 2: Ablation study results showing the impact of removing or altering components of the proposed HFFN model, evaluated by ROUGE-L score.

Model Variant	ROUGE-L	Δ vs. Full HFFN
HFFN (Full Model)	0.94	-
- without Syntactic Features	0.85	-0.09
- without Lexical Features	0.87	-0.07
- with Concatenation Fusion	0.90	-0.04

- **Impact of Syntactic Features:** The removal of the syntactic feature branch (AST-based features) resulted in the largest performance drop (-9.6%). This confirms that syntactic structure provides the greatest individual impact on model performance, forming the essential backbone of an accurate code representation.
- **Impact of Lexical Features:** Removing the lexical feature stream (sequence-based tokens) also caused a significant degradation in performance (-7.4%). This demonstrates that surface-level tokens, including meaningful API and variable names, provide indispensable semantic signals.
- **Impact of Fusion Mechanism:** Replacing the gated fusion mechanism with simple feature concatenation led to a notable decline (-4.3%). This validates the design hypothesis that a learned, adaptive fusion strategy is superior, as it allows the model to dynamically weight the contribution of each feature stream for a given input[16].

Discussion: The ablation study confirms that the components of HFFN provide complementary contributions. The architecture successfully leverages the strengths of both feature types, with the gated fusion mechanism acting as a critical component for intelligently combining these information streams. The significant drop from removing either modality underscores the necessity of a hybrid approach[13].

4.4 Qualitative Analysis

A manual inspection of generated summaries revealed consistent strengths and weaknesses. For a function implementing a quicksort algorithm, baseline models often produced vague summaries like “Sorts an array”. In

contrast, HFFN generated more precise and informative summaries such as “Sorts an integer array in ascending order using the quicksort algorithm with a median-of-three pivot selection”.

The model’s limitations were apparent on functions with extreme complexity or those utilizing rare external libraries, where summaries would occasionally miss a minor detail or over-generalize. These cases represent the current boundary of the model’s capabilities.

4.5 Limitations

While HFFN demonstrates strong performance, its effectiveness is contingent on the availability of a well-formed Abstract Syntax Tree (AST), rendering it unsuitable for processing syntactically incorrect code snippets. Furthermore, the computational overhead associated with AST parsing and processing the three specialized encoders presents a trade-off between performance and efficiency that may be relevant for real-time applications with limited resources.

While the proposed HFFN model achieves high performance on the CanonCode Corpus, it is important to note that the dataset primarily consists of textbook examples and tutorial code. Such code tends to be well-structured, pedagogically motivated, and accompanied by standardized summaries. This may contribute to the high ROUGE-L and BERTScore values observed. Future work will involve testing HFFN on more diverse and complex real-world codebases to further validate its generalizability.

5 Conclusion and Future Work

5.1 Conclusion

This study presented the Hybrid Feature Fusion Network (HFFN), a novel architecture for source code representation that explicitly models and integrates syntactic and lexical features through a gated attention mechanism. The comprehensive evaluation across five established metrics—ROUGE-L, BLEU, Exact Match, BERTScore, and CodeBLEU—demonstrates the efficacy of this approach. The key conclusions are as follows:

- **Superior Performance Across Metrics:** HFFN established a new state-of-the-art, outperforming all strong baselines, including CodeBERT, across every evaluation metric. The model achieved a ROUGE-L score of 0.94 and a BERTScore of 0.93, indicating its generated summaries are both semantically faithful and fluent. The significant improvement in the strict Exact Match metric (a 50% relative increase) confirms that HFFN produces perfectly accurate summaries far more frequently than previous approaches.
- **Validation of Hybrid Methodology:** The results provide convergent evidence that a hybrid approach is paramount for code understanding. The notable gains in CodeBLEU, a metric specifically designed to evaluate code by incorporating syntactic matching, directly validate our core hypothesis: explicitly modeling the structural hierarchy of code is a powerful and necessary inductive bias. The strong performance on BLEU further demonstrates the model’s syntactic precision in using correct terminology.
- **Effective Fusion Mechanism:** The ablation study conclusively proved that both syntactic and lexical features provide significant, complementary contributions to the model’s performance. Furthermore, the gated fusion mechanism was shown to be a critical component, significantly outperforming a naive feature concatenation strategy. This allows the model to dynamically and intelligently weight the contribution of each feature stream.

The proposed HFFN provides a robust, interpretable, and high-performing framework for source code representation. By successfully fusing disparate code features into a cohesive representation, it offers a powerful foundation for a new generation of AI-assisted software engineering tools. This study introduced a Hybrid Feature Fusion Network (HFFN) for code summarization that combines lexical, syntactic, and semantic information. Experimental results showed that this multi-view approach consistently surpasses methods relying on a single type of feature representation, including those built on large-scale pre-trained models. The findings highlight the value of integrating complementary perspectives of source code to produce clear, high-quality summaries.

5.2 Future Work

While this study presents a significant advancement, several promising directions for future research emerge from the findings and limitations:

- **Extension to Other Programming Languages and Paradigms:** The current evaluation was conducted on C programs. A critical next step is to evaluate HFFN's generalizability across a wider range of programming languages, particularly those with different paradigms (e.g., functional languages like Haskell, or scripting languages like Python). This will test the universality of the proposed feature fusion approach.
- **Integration of Execution Dynamics and Real-World Context:** The current model operates on static code features. Future work will focus on incorporating dynamic features, such as runtime traces and execution patterns, to create a more holistic code representation. Furthermore, integrating broader contextual information from the codebase, such as project-specific APIs and architectural patterns, could significantly enhance the model's practical utility.
- **Optimization for Efficiency and Real-Time Use:** The computational overhead associated with AST processing presents a constraint for real-time applications. Future work will explore model distillation techniques, efficient AST parsers, and architectural optimizations to make the HFFN approach viable for integration into Interactive Development Environments (IDEs) for tasks like real-time code completion and documentation generation.
- **Advanced Fusion and Explainability Techniques:** We will investigate more sophisticated fusion mechanisms, potentially inspired by neuroscientific models of information integration. Concurrently, developing advanced visualization tools based on the model's attention weights will enhance explainability, allowing developers to understand why a particular summary or prediction was generated, thereby increasing trust and usability.
- We also plan to evaluate HFFN on established public benchmarks such as CodeXGLUE to further validate its generalizability across different programming languages and codebases.

References

- [1] Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. Transformer-based models for code summarization: A systematic review. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4998–5007, 2020.
- [2] Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. Unified pre-training for program understanding and generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2668, 2021.
- [3] Miltiadis Allamanis, Earl T Barr, Premkumar Devanbu, and Charles Sutton. A survey of machine learning for big code and naturalness. *ACM Computing Surveys*, 51(4):81, 2018.
- [4] Miltiadis Allamanis, Marc Brockschmidt, and Mahmoud Khademi. Learning to represent programs with graphs. In *International Conference on Learning Representations (ICLR)*, 2018. Preprint available at arXiv.
- [5] Miltiadis Allamanis, Hao Peng, and Charles Sutton. A convolutional attention network for extreme summarization of source code. In *International Conference on Machine Learning*, pages 2091–2100. PMLR, 2016.
- [6] Shruthi D., H. K. Chethan, and Agughasi Victor Ikechukwu. A novel deep learning-based technique for automatic source code summarization. In Milan Tuba, Shyam Akashe, and Amit Joshi, editors, *ICT Systems and Sustainability*, pages 525–537, Cham, 2026. Springer Nature Switzerland.
- [7] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. Codebert: A pre-trained model for programming and natural languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1536–1547, 2020.
- [8] Daniel Gros, Hrant Sezhiyan, Premkumar Devanbu, and Zhou Yu. Code to comment translation: A comparative study on model effectiveness and errors. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 837–848, 2021.
- [9] Sonia Haiduc, Jairo Aponte, and Andrian Marcus. Supporting program comprehension with source code summarization. In *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering—Volume 2*, pages 223–226, 2010.
- [10] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

- [11] Xing Hu, Ge Li, Xin Xia, David Lo, and Zhi Jin. Deep code comment generation. In *Proceedings of the 26th Conference on Program Comprehension*, pages 200–210, 2018.
- [12] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. Codesearchnet challenge: Evaluating the state of semantic code search. *arXiv preprint*, 2019.
- [13] Agughasi Victor Ikechukwu. Leveraging transfer learning for efficient diagnosis of copd using cxr images and explainable ai techniques. Department of Computer Science and Engineering.
- [14] Agughasi Victor Ikechukwu. The superiority of fine-tuning over full-training for the efficient diagnosis of copd from cxr images. Department of Computer Science and Engineering.
- [15] Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. Summarizing source code using a neural attention model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2073–2083, 2016.
- [16] Reem Jalloul, Chethan Hasigala Krishnappa, Victor Ikechukwu Agughasi, and Ramez Alkhatib. Enhancing early breast cancer detection with infrared thermography: A comparative evaluation of deep learning and machine learning models. *Technologies*, 13(1), 2025.
- [17] Alexander LeClair, Siyuan Jiang, and Collin McMillan. A neural model for generating natural language summaries of program subroutines. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*, pages 795–806. IEEE, 2020.
- [18] Paul W McBurney and Collin McMillan. Automatic documentation generation via source code summarization of method context. In *Proceedings of the 22nd International Conference on Program Comprehension*, pages 279–290, 2014.
- [19] Paul W. McBurney and Collin McMillan. Automatic documentation generation via source code summarization of method context. In *Proceedings of the 22nd International Conference on Program Comprehension, ICPC 2014*, page 279â290, New York, NY, USA, 2014. Association for Computing Machinery.
- [20] Prathibha S., Madhusudhan K. N., and Agughasi Victor Ikechukwu. Firefly-based segmentation and residual deep learning for multi-class diabetic retinopathy detection. Affiliated institutions: BMS College of Engineering, Bengaluru and Maharaja Institute of Technology Mysore.
- [21] D. Shruthi, H. K. Chethan, and Agughasi Victor Ikechukwu. Evaluating the feasibility of a pre-processing framework for enhanced text information extraction from source codes. In Jagdish Chand Bansal, Snehanu Saha, Carlos A. Coello Coello, and Hemant Rathore, editors, *Advances in Data-Driven Computing and Intelligent Systems*, pages 335–350, Singapore, 2025. Springer Nature Singapore.
- [22] D Shruthi, H.K. Chethan, and V.I. Agughasi. Effective approach for fine-tuning pre-trained models for the extraction of texts from source codes. In *ITM Web of Conferences*, volume 65, page 03004. EDP Sciences, 2024.
- [23] Kai Sheng Tai, Richard Socher, and Christopher D Manning. Improved semantic representations from tree-structured long short-term memory networks. *arXiv preprint*, 2015.
- [24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- [25] Anthony J Viera and Joanne M Garrett. Understanding interobserver agreement: The kappa statistic. *Family Medicine*, 37(5):360–363, 2005.
- [26] Yao Wan, Zhou Zhao, Min Yang, Guandong Xu, Haochao Ying, Jian Wu, and Philip S. Yu. Improving automatic source code summarization via deep reinforcement learning. *arXiv preprint*, 2018.
- [27] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8696–8708, 2021.
- [28] Bolin Wei, Ge Li, Xin Xia, Zhiyuan Fu, and Zhi Jin. Code generation as a dual task of code summarization. In *Advances in Neural Information Processing Systems*, pages 6559–6569, 2019.
- [29] Yuxiang Zhou, Shi Liu, Jingkai Siow, Xiaoning Du, and Yang Liu. Code summarization with structure-induced transformer. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*, pages 1078–1088. IEEE, 2019.