
AUTOMATIC REPAIR AND TYPE BINDING OF UNDECLARED VARIABLES USING NEURAL NETWORKS

A PREPRINT

Venkatesh Theru Mohan
 Department of Computer Science
 Iowa State University
 venkytm@iastate.edu

Ali Jannesari
 Department of Computer Science
 Iowa State University
 jannesar@iastate.edu

July 16, 2019

ABSTRACT

Deep learning had been used in program analysis for the prediction of hidden software defects using software defect datasets, security vulnerabilities using generative adversarial networks as well as identifying syntax errors by learning a trained neural machine translation on program codes. However, all these approaches either require defect datasets or bug-free source codes that are executable for training the deep learning model. Our neural network model is neither trained with any defect datasets nor bug-free programming source codes, instead it is trained using structural semantic details of Abstract Syntax Tree (AST) where each node represents a construct appearing in the source code. This model is implemented to fix one of the most common semantic errors, such as undeclared variable errors as well as infer their type information before program compilation. By this approach, the model has achieved in correctly locating and identifying 81% of the programs on prutor dataset of 1059 programs with only undeclared variable errors and also inferring their types correctly in 80% of the programs.

Keywords Automatic Bug Repair, Undeclared Variables, Neural Networks, Type Binding

1 Introduction

In the recent decades, there had been some significant contributions in the automation of program analysis tasks. Neural networks had been frequently used either for detection or generation tasks in programming language processing similar to natural language processing. It had been employed in detecting software defects, as well as prediction of errors by using software metrics [17].

There have been some state-of-art performances achieved by generative adversarial networks and neural machine translation systems on language translation tasks in Natural Language Processing (NLP) that led to the deployment of these systems on error correction tasks in the programming source codes. Neural machine Translation (NMT) in [14] builds and trains a single neural network model using a labelled set of paired examples that results in translation from the input directly. They are end-to-end in the sense that they have the ability to learn the source text directly by mapping them to the corresponding target text and uses an encoder-decoder(usually made up of Recurrent Neural Networks (RNN) units) approach for applying the transformations where encoder consists of sequences of source text and output consists of sequences of target text. Generative Adversarial Networks (GAN) trains a neural network model to predict security vulnerabilities in [13] without the necessity of being a paired set of source domain containing buggy codes and target domain consisting of bug-free codes, instead being a bijection mapping. Generally, in this labelled set of paired examples or unpaired examples, the neural network model is trained on the set of positive examples where the mapping takes place between target sequences, made up of positive examples that are bug-free and compiling without any errors and source sequences consisting of negative examples that contains bugs within the code making the compilation to fail. In sequence-to-sequence learning systems such as NMT, it is flexible to train and learn the model with a labelled set of paired examples, but consider the scenario where there are no paired

examples, further to put it in a simpler context, consider a real-case scenario where some common syntax or semantic errors are being committed by freshers or novice programmers working in a software industry or student submissions of programming assignments in coding competitions and there are no positive or bug-free reference examples, then learning becomes difficult for neural network approaches and neural machine translation models.

In our model, there are no positive examples to train and focus is only on the negative buggy examples specifically, the common semantic error caused due to undeclared variables which is often committed and unperceived by the novice programmers. The model is instead trained based on the structural and semantic elements, that is the non-terminal and terminal nodes of the programming source codes captured from the abstract syntax tree representation. The type of the undeclared variables is also inferred by performing type binding using the semantic elements of AST representation that provides about the type information of the variables thereby saving the compiler's time in performing the type binding of those undeclared variables. The comprehensive information of the ASTs, the motive of Long Short-Term Memory (LSTM), detailed view of the training approach and implementation, the generation approach and different scenarios where undeclared variables will be caused and possible cases of type inference of them will be discussed in the upcoming sections of the paper.

2 Related Work

AST are the static intermediate model of a programming source code as discussed in [1] where the compiler's analytic front-end parses it, constructs the AST model eventually passing it to the compiler's synthetic back-end to produce an assembly code for a specific machine and also used for program analysis/transformation. Often low-level concrete syntax tree representations have a complex structure and it is difficult to characterize the semantics especially in the poorly understood domains, so in [2], describes a transformation process to get a good abstract syntax representation from low-level concrete specification where modern tools rely on its ability to analyze, simulate and synthesize programs easily in language processing.

In the past recent decades, deep learning had achieved various scales in text domain such as natural language especially in neural machine translation tasks and used it successfully even for programming language processing tasks as well. Some of the recent works that had achieved empirical success on neural machine translation include dialogue response generation, summarization and text generation tasks as explained in detail in [3], also NMT tasks are very much useful in translation from one language to another like in [4] where NMT encoder-decoder model had been implemented using word embedding along with byte-pair encoding method to develop an efficient translation mechanism from English-Tamil. The performance of such text generation tasks such as NMT is enhanced with attention mechanism in [5] to automatically search for a context of a source text that is relevant for prediction of target words and an ensemble model of global and local attention mechanism of [6] improving the performance.

Natural language generation having a discrete output space had also been implemented by generative adversarial networks (introduced by [7]) on generating sentences from context-free grammars and probabilistic grammars as shown in [8]. The quality of the text generation had been improvised by conditioning information on generated samples to produce realistic natural looking sentences compared to maximum likelihood training procedures. Text generation has also shown some encouraging results in [9] in which the discriminator of the GAN leaks its own high-level features to the generator at each of the generator steps thus making the scalar guiding signal available continuously throughout the generative process.

In [10], prediction of next tokens in the source code is implemented using LSTM neural networks where the model trains and learns associated subsequent nodes for code completion, given a partial AST containing left subset of nodes or semantic features with respect to a subtree. The efficiency of AST in extracting tokens and comparing the source codes based on them and the use of deep learning in classifying duplicate/clone codes can be helpful in software code maintenance as seen in [11] where maintaining duplicate codes for reuse in order to improve productivity of programming becomes a burden when there are inconsistencies caused due to bug fixes and program modifications in the original code at multiple locations.

The significance of AST representation of source code can further be noted in [12] where LSTM neural networks are leveraged to capture the long contextual relationships between semantic features to identify related code elements in order to predict the software vulnerabilities which causes a security threat or makes the program buggy. GAN approach is used in [13] for repairing vulnerabilities in source codes without any paired examples or bijections by mapping from non-buggy source domain to buggy target domain and training the discriminator using the loss that is generated between real examples and NMT-generated outputs from generator of desired output.

```

1  #include <stdio.h>
2  int main()
3  {
4      int i, max, j, n, m, y;
5      scanf("%d %d", &n, &m);
6      for(i = 1; i <= n; i++){
7          s = 0;
8          for(j = 1; j <= m; j++){
9              scanf("%d", &y);
10             s = s + j;
11         }
12         if(max < s)
13             max = s;
14     }
15     return 0;
16 }

```

Figure 1: An example illustrating the undeclared variable "s" that is frequently used in the program is caught by the compiler

Syntax errors poses a threat as it fails the compilation and some recent techniques such as [14] where a RNN model is learnt on syntactically correct, executing student programming course submissions to model all valid sequences of token and use a prefix token sequence which is from the beginning of the program till the error location and is used to predict the following sequence that are able to automate the repair of errors present in corresponding locations of code. Sequence-to-sequence NMT with attention mechanism is learned iteratively in repairing syntax errors in [15] using the tokenized vector representation of the program and is used to predict the erroneous program locations and the fixing statement without using any external compiler tools or any AST representation. A real-time feedback is given to the students enrolled in beginner level programming assignments in [16] of the compile-time syntax errors that are made by using RNN to predict the target lines from syntactically correct submissions given the source error lines from wrong submissions and a abstract version of top ranked suggestion of error fix is presented as feedback.

3 Approach

This section covers some examples of undeclared variables and also the importance of semantic analysis to determine the type information of those undeclared variables. We introduce the Abstract Syntax Trees (AST) that is used as the input, deployment of the LSTM RNN for training the deep learning model, semantic analysis determining the types of undeclared variables, and the generation approach by performing the serialization/deserialization of the AST in order to get back the clean and bug-free source code.

3.1 Motivating Examples

The most frequent semantic error that goes unnoticed by novice programmers is the undeclared variables. The cause of this error is due to the variables being undeclared or another common cause will be usually through spelling mistakes which makes it the first occurrence in the program. The main challenge lies in the fact in determining whether the variable is an identifier, arrays, pointers or pointers-to-pointers and also in concluding about the type of the variable if it is an integer, float, character, double, long int and so on. The C99 standard¹, removes the implicit integer rule that states a variable declared without an explicit data type is assumed to be integer which was previously defined in C89 standard². Therefore, there is a need for determining the variables that are undeclared along with its type, else the compiler will throw an error.

3.2 Abstract Syntax Trees

Generally, any programming language whether statically or dynamically typed follows an unambiguous context-free grammar language where there exists not more than one leftmost derivations or rightmost derivations of non-terminals, or more precisely there is always a unique parse tree for each string of the language generated by it. Parse trees are

¹www.open-std.org/jtc1/sc22/wg14/www/docs/n1256.pdf

²<https://www.pdf-archive.com/2014/10/02/ansi-iso-9899-1990-1/ansi-iso-9899-1990-1.pdf>

```

1  #include <stdio.h>
2  int main ()
3  {
4      int n,m;
5      int i , j ;
6      int a [ 20 ];
7      int sum=0;
8      scanf ( "%d%d",&n,&m) ;
9      for( i =0; i <n; i ++){
10         for( j =0; j <m; j ++){
11             scanf ( "%d",&a [ j ] ) ;
12             sum=sum+a [ j ] ;
13         }
14         printf ( "%d\n",sum) ;
15         i ++;
16         J ++;
17     }
18     return 0;
19 }

```

Figure 2: An example of compiler error caused by an undeclared variable "J" that is used once in the program

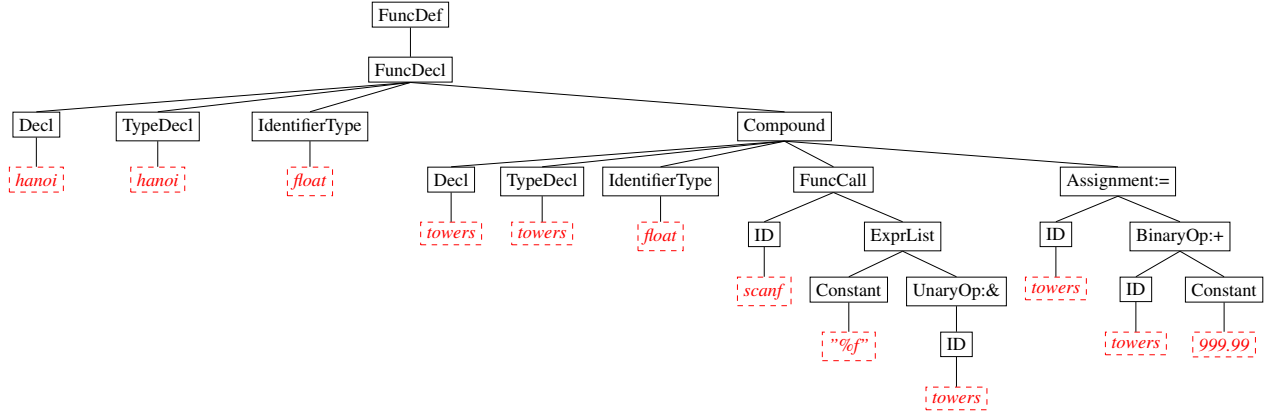


Figure 3: An example illustration of Abstract Syntax Trees(AST) of C program with non-terminals at the root as well as internal nodes and terminals at the leaf nodes of the tree.

formed from a concrete context-free grammar and are not suitable for performing syntax or semantic analysis due to their complex representation. Nevertheless, the Abstract Syntax Tree are the derivation trees following an abstract grammar is used as an input for syntax/semantic analysis at compile-time. It is a rooted tree representation of an abstract syntactical structure of a programming language construct where the non-terminals represents non-leaf nodes and the terminals forms the children nodes. Some of the non-terminals in C language are Decl, IdentifierType, For, TypeDecl, ExprList, FuncDef, ArrayDecl etc. Terminal nodes includes any string literals, numerical literals, variable names, operators, keywords etc. Figure 3 shows an example AST representation where a node enclosed in a rectangle box and black-colored depicts the non-terminal nodes (eg.,FuncDef), node that are enclosed in dashed box and red-colored depicts the terminal nodes or tokens of the source code.

3.3 Model

This subsection covers the basics of LSTM-RNN and the prediction model that is subsequently used for generation approach.

Long Short Term Memory(LSTM) recurrent neural networks have special memory units in the form of self-loops to produce paths so that information can be maintained for longer durations of time. LSTM is preferred over vanilla RNN due to the fact that the former tends to avoid vanishing or exploding gradient problem that occurs when trying to learn long term dependencies and store it in memory cells during backpropagation. This occurs when many deep layers with specific activation functions like sigmoid are used for training, it smoothens a region of input space

into an output space between 0 and 1, then even a high change in input region effects almost negligible change in output region, thereby making the gradients/error signals of a long-term interaction becomes vanishingly very small. Further, the vanilla RNNs are affected by information morphology problem in which information contained at a prior state is lost due to non-linearities between the input and output space. LSTM avoids this problem by ensuring a constant unit activation function and uses gates to control the information flow between the memory cell and the outside layers without any inference. LSTM uses three gates namely forget gate, input gate and output gate layers. The forget gate layer decides the information that is needed to be stored or erased from the LSTM cell state where the decision is made by the sigmoid layer outputting a number between 0 to 1.

$$f_i^{(t)} = \sigma(b_i^f + \sum_j U_{ij}^f x_j^{(t)} + \sum_j W_{ij}^f h_j^{(t-1)}) \quad (1)$$

where $x^{(t)}$ is the input vector at current timestep t and $h^{(t)}$ is the current hidden layer vector at timestep t , and b^f , U^f , W^f are bias units, input weights and recurrent weights of forget gate units $f_i^{(t)}$.

The input gate layer controls the flow of new information that is being stored in the LSTM cell state $s_i^{(t)}$ conditioned with a self-loop weight $f_i^{(t)}$.

$$g_i^{(t)} = \sigma(b_i^g + \sum_j U_{ij}^g x_j^{(t)} + \sum_j W_{ij}^g h_j^{(t-1)}) \quad (2)$$

where $x^{(t)}$ is input vector at current timestep t and $h^{(t)}$ is current hidden layer vector at timestep t , and b^g , U^g , W^g are bias units, input weights and recurrent weights of input gate units $g_i^{(t)}$.

$$s_i^{(t)} = f_i^{(t)} s_i^{(t-1)} + g_i^{(t)} \sigma(b_i + \sum_j U_{ij} x_j^{(t)} + \sum_j W_{ij} h_j^{(t-1)}) \quad (3)$$

where the parameters W, U and b represents the recurrent weights, input weights and bias units present in a LSTM cell. The output gate layer in the memory cell decides the pieces of information that is going to be output by the LSTM cell state. This is done by passing the cell state through a tanh layer and eventually multiplying by the sigmoid of the output gate.

$$q_i^{(t)} = \sigma(b_i^o + \sum_j U_{ij}^o x_j^{(t)} + \sum_j W_{ij}^o h_j^{(t-1)}) \quad (4)$$

where b^o, U^o, W^o are the parametric units of the output gate $q_i^{(t)}$ that represents bias units, input weights and recurrent weights. The output hidden state $h_i^{(t)}$ is obtained from output gate $q_i^{(t)}$ as follows:

$$h_i^{(t)} = \tanh(s_i^{(t)}) q_i^{(t)} \quad (5)$$

In our prediction model, we use non-terminals and terminals obtained from the AST as the input. The main ideology behind our model is to specify a declaration for any identifier that is used throughout a C program. Here, we assume identifier in our context that excludes keywords and only includes alphanumeric variables. This in turn solves the complex problem of automatically fixing the undeclared variables.

3.3.1 Declaration Classification and Prediction

ID is the non-terminal node of the AST that represents an identifier excluding keywords as mentioned above. Decl is the non-terminal node of the AST representation that is entitled to represent the declaration of the identifier where its terminal node is the corresponding identifier itself. Similarly, TypeDecl and IdentifierType are the semantic elements that is used to represent the type specifier information of the identifier where identifier and its type are its corresponding terminals respectively.

For the classification purpose, the pair of non-terminal node ID and any terminal alphanumeric identifiers usually variables is augmented along with pairs of Decl and the respective identifier, TypeDecl and the identifier, IdentifierType and a generalized "type" referring to the corresponding types of those identifier variables so that they can be used for backsubstitution which will be explained later in the generation approach. After classification, the LSTM model is used to predict the Decl, TypeDecl and IdentifierType information for any alphanumeric identifier variable occurring with its corresponding non-terminal node ID. The classification model is sequential where the embedding layer, LSTM layer and softmax layers are stacked on top of each other sequentially as shown in Figure 4 are detailed below in this section.

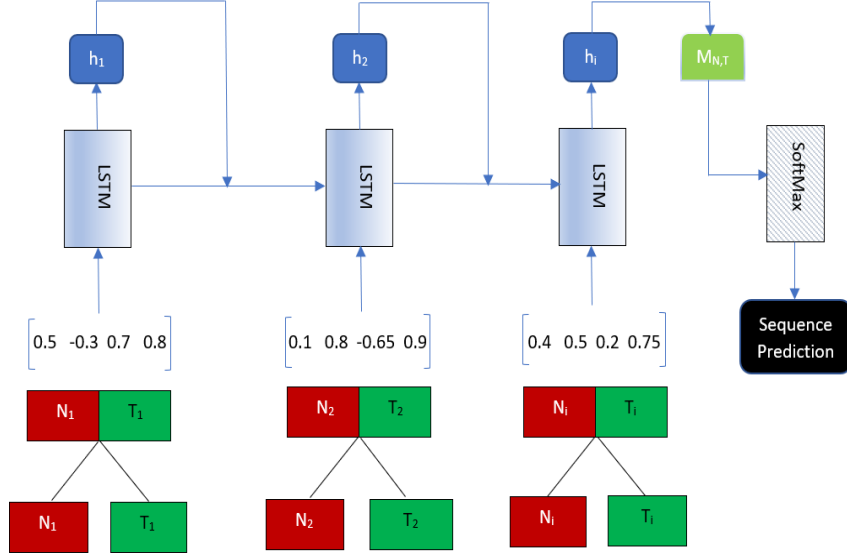


Figure 4: Illustration of approach showing concatenation of non-terminal and terminal node embeddings extracted from AST being used as the inputs to LSTM model for the sequence classification and prediction.

3.3.2 Embedding Layer of Non-Terminal and Terminal Nodes

In our model, we do not use any pre-trained embeddings, instead it is trained simultaneously with the model. The sequences of input tokens are a combination of non-terminal and terminal nodes where in our model, we consider only 4 non-terminals ID, Decl, TypeDecl, IdentifierType and all the alphanumeric identifier variables as the terminals. These input tokens are formed by the concatenation of individual string encodings of non-terminal and terminal node vocabularies (that is discussed in evaluation section), subsequently embedding is computed on the integer encodings(converted from string) for performing the training of the model. The embeddings are computed as follows:

$$E_i = A * concat(N_i T_i) \quad (6)$$

where A is $K \times V_{N,T}$ matrix where K is the embedding vector size and $V_{N,T}$ is the vocabulary size formed by the concatenated encodings of non-terminal and terminal nodes.

3.3.3 LSTM Layer

The sequences of embedded tokens are passed on to the LSTM layer containing LSTM memory cells where each cell state stores information of the previous state and are controlled by the forget gate, input gate and output gate layers as mentioned above in equations (1),(2),(3),(4). Each LSTM cell state takes inputs from its previous LSTM cell hidden state h_{i-1} , state information s_i as well as the input tokens and outputs the hidden state h_i of LSTM cell as in equation (5) where the LSTM layers can be seen from Figure 4.

3.3.4 Dense Softmax Activation Layer

The last LSTM memory cell state's output hidden state of the LSTM layer is passed to the softmax activation layer to predict the sequences of non-terminals Decl, TypeDecl and IdentifierType given the non-terminal ID. The predicted output sequences at timestep t (or fixed input sequence length) are represented by $\hat{y}(t)$ and is formulated as:

$$\hat{y}(t) = softmax(b_{N,T} + M_{N,T} h_{(t)}) \quad (7)$$

where $b_{N,T}$ is the bias unit of softmax layer with size of $V_{N,T}$ dimensional vector, $M_{N,T}$ is a weight matrix of size $K \times V_{N,T}$ and $h_{(t)}$ is the hidden state of the LSTM cells at each timestep t .

```
{1: 'Decl', 2: 'IdentifierType', 3: 'Constant', 4: 'FuncDef', 5: 'FuncDecl',
6: 'FuncCall', 7: 'ParamList', 8: 'Typename', 9: 'TypeDecl', 10: 'Compound',
11: 'Assignment', 12: 'ID', 13: 'UnaryOp', 14: 'BinaryOp', 15: 'Cast',
16: 'ExprList', 17: 'Pragma', 18: 'For', 19: 'While', 20: 'DoWhile', 21: 'If',
22: 'Switch', 23: 'Return', 24: 'Print', 25: 'DeclList', 26: 'ArrayRef',
27: 'Cast', 28: 'TypeDef', 29: 'PtrDecl', 30: 'Struct', 31: 'ArrayDecl',
32: 'InitList', 33: 'Break', 34: 'Case', 35: 'CompoundLiteral', 36: 'Continue',
37: 'Default', 38: 'EmptyStatement', 39: 'Enum', 40: 'Enumerator',
41: 'EnumeratorList', 42: 'Goto', 43: 'Label', 44: 'NamedInitializer',
45: 'StructRef', 46: 'TernaryOp', 47: 'Union'}
```

Figure 5: The vocabulary of all non-terminal nodes of AST for C programs

4 Evaluation and Results

4.1 Dataset

The dataset that used in this approach is prutor³ which is a database that has student coding submissions for university programming assignments. It contains a set of 53478 C programs out of which there are 6978 erroneous programs which contains multiple and single line syntax as well as semantic errors. Out of 6978 programs, 1059 programs contains only undeclared variable errors which is the main focus of our evaluation.

4.2 Preprocessing and Training Details

We use pyparser⁴, that acts as a front-end of the C compiler to parse source code of C language in python. AST are obtained as an output for the source code after the parsing stage and are stored in text files.

The source code is preprocessed in the form of tokens that represents the terminal nodes of its corresponding AST. Since the set of tokens are uncontinuous, discrete and in its textual form, it must be encoded into sequences of numerical vectors to be used for training the model. Additionally, there are 47 fixed set of non-terminals in C language that are encoded as in Figure 5. The terminals can be keywords, strings, data types, integers or floating point numbers. The terminals of the above mentioned categories are encoded separately in a specific range of numbers. Now, the data for training is prepared by concatenating the encodings of non-terminals and terminals together. The individual encodings in the form of integers are converted to strings initially and after concatenation, they are converted back to integers. For example, the non-terminal IdentifierType is encoded as 9 in the dictionary of non-terminals and converted to '9', if the terminals are data types like int, float, long etc. then they are encoded as 111111 referring to a generalized 'type' and converted to '111111'. Therefore, the concatenation of the non-terminals and terminals are mapped accordingly and stored in separate vocabulary. This vocabulary set is used in performing the training of the model. One-hot encoding approach is used to perform categorical multi-class classification to represent the elements of vocabulary as vectors with each of them of vocabulary size containing 1 at the corresponding index and rest of them are 0.

Training Details:

The training experiment is performed by using embedding dimension of size 512 and two LSTM layers are used each with number of hidden units as 512 and a dropout of 0.5. The input sequence length is 1, batch size is 3, vocabulary size is 583 and the total number of sequences is 2319. The dense layer is used for forming a fully connected layer in which each of the input layer nodes is connected with every hidden and output layer nodes. The activation function used in our model is softmax function because of its efficiency in dealing with multi-class classification problems compared to sigmoid and ReLU due to the fact that the outputs of the softmax is a categorical probability distribution summing to 1 and lying between 0 to 1. The total number of units of Dense layer is equal to the vocabulary size. The vocabulary formed from concatenation is split into training and testing data where test data size is 0.2 and the split rule percentage used is 80/20. The loss function used is categorical cross-entropy function. The optimizer used is RMSprop with a learning rate of 0.01 as it is better in handling non-local extremum points and has a constant initial global learning rate compared to optimizers such as Stochastic gradient descent optimizer.

³<https://www.cse.iitk.ac.in/users/karkare/prutor/prutor-deepfix-09-12-2017.zip>

⁴<https://pypi.org/project/pyparser/>

```

{
  "_nodetype": "Decl",
  "bitsize": null,
  "funspec": [],
  "init": null,
  "name": "j",
  "quals": [],
  "storage": [],
  "type": {
    "_nodetype": "TypeDecl",
    "declname": "j",
    "quals": [],
    "type": {
      "_nodetype": "IdentifierType",
      "names": [
        "int"
      ]
    }
  }
}

```

Figure 6: Example demo of JSON object containing Decl, TypeDecl and IdentifierType nodes.

4.3 Generation Approach

During generation, one-hot encoding is used to represent the new unseen sequences of the nonterminal ID and a terminal variable as a categorical distribution and the output class is classified as the sequences of non-terminals Decl, TypeDecl, IdentifierType along with the corresponding terminal variables.

4.3.1 AST Transformation

The program fix approach is carried out through an AST transformation performed by augmenting the predicted output sequences in each of the program’s AST syntactical structure for the terminal variables associated with its corresponding non-terminal node ID. This augmentation is carried out on each of the source code by performing a check on declaration of the variables in the vocabulary set of concatenated encodings of non-terminal and terminal nodes as created previously using the predicted output sequence of the non-terminal Decl and the associated terminal variable. If any of the predicted output sequences does not match with the declared variables present in a source code, then the output sequence Decl containing the particular terminal variable is augmented to the original AST structure of the code through serialization and deserialization that will be described below.

4.3.2 Serialization and Deserialization

Serialization is implemented using pycparser by transforming data structures such as nodes of python Node object by traversing the AST (the nodes being obtained by parsing the source code from pycparser) recursively into a dictionary object representation which is then subsequently serialized into a JSON object that is understood by Pycparser in deserializing it back to a dictionary object and thereby consequently back to AST Node objects. An example JSON object representation of a AST node object is shown in Figure 6 where the _nodetype key refers to the different types of non-terminal nodes such as Decl, ArrayDecl, TypeDecl, IdentifierType. The TypeDecl and ArrayDecl are the child nodes of Decl and IdentifierType is the child node of the intermediate non-terminal TypeDecl node. The key name refers to the terminal variables, type refers to the datatypes and coord refers to the Coord node of the AST indicating the location of the object in the corresponding source code.

The serialization and deserialization is carried out when there is an AST transformation performed as mentioned above, so that the transformation is taking place consistently without disintegrating the program thereby ensuring maximal quality of the source code.

4.3.3 Pre-Compile Time Type Binding and Analysis Results

After determining and performing the augmentation of the undeclared terminal variables, the type of those undeclared variables is determined before compiling the program and is augmented to the original AST structure. The type binding is performed by determining the type of a lvalue from its rvalue in an assignment statement or finding the type of a undeclared variable from its neighboring variables in an expression of a program.

The type of the undeclared variables is determined and drawn from the following cases:

Case 1: When the rvalue is a constant, where the non-terminal is Constant and is of integer type, then the lvalue whose non-terminal is ID, is assigned integer as seen from the Figure 7, in the assignment statement **i=0** on left side of figure, "i" is undeclared and is assigned integer.

<pre> 1 int main() 2 { 3 int k,n,x,a[100]; 4 scanf("%d",&k); 5 scanf("%d",&n); 6 for(i=0;i<n;i++) 7 scanf("%d",&a[i]); 8 return 0; 9 }</pre>	<pre> 1 int main() 2 { 3 int i; 4 int k; 5 int n; 6 int x; 7 int a[100]; 8 scanf("%d",&k); 9 scanf("%d",&n); 10 for (i = 0; i < n; i++) 11 scanf("%d",&a[i]); 12 return 0; 13 }</pre>
---	---

Figure 7: Case 1 demonstrating location of error in the for loop statement involving undeclared identifier "i" and the fix of it

Case 2: In this case, if the rvalue is an identifier that refers to an array element whose non-terminal is ID and its non-terminal parent is ArrayRef, then the lvalue terminal variable with non-terminal ID is assigned to the respective type of rvalue element. We can see in Figure 8 where "b" is undeclared and is assigned to integer type from array "n[1000]" in the statement **b=n[i]**.

<pre> 1 int main() 2 { 3 int n[1000],a[500],nm,i,j,ln, 4 flag=0; 5 scanf("%d\n",&ln); 6 scanf("%d\n",&nm); 7 for(i=0;i<500;i++) 8 { 9 a[i]=0; 10 } 11 for(i=0;i<nm;i++) 12 { 13 scanf("%d",&nm); 14 c=n[i]; 15 } 16 return 0; 17 }</pre>	<pre> 1 int main() 2 { 3 int c; 4 int n[1000]; 5 int a[500]; 6 int nm; 7 int i; 8 int j; 9 int ln; 10 int flag = 0; 11 scanf("%d\n",&ln); 12 scanf("%d\n",&nm); 13 for (i = 0; i < 500; i++) 14 { 15 a[i] = 0; 16 } 17 for (i = 0; i < nm; i++) 18 { 19 scanf("%d",&nm); 20 c = n[i]; 21 } 22 return 0; 23 }</pre>
---	---

Figure 8: Case 2 indicating the error in assignment statement between variable and array identifier

Case 3: If the non-terminal of rvalue and non-terminal of lvalue are the children nodes of the non-terminal BinaryOp, then the type of rvalue is assigned as the type of lvalue. In the Figure 9 on the left side of the figure, in the conditional expression statement $count > max$, BinaryOp is ">", the lvalue "count" is undeclared with its non-terminal being ID and the rvalue is max with its non-terminal being ID.

<pre> 1 int main() 2 { 3 int n, i, j, max; 4 int a[20]; 5 for(i=0; i<n; i++) 6 { 7 for(j=i; j<n; j++) 8 { 9 if(a[i]<a[j]) 10 { 11 count=count+1; 12 } 13 } 14 if(count>max){max=count;} 15 } 16 printf("%d",max); 17 return 0; 18 } </pre>	<pre> 1 int main() 2 { 3 int count; 4 int n; 5 int i; 6 int j; 7 int max; 8 int a[20]; 9 for (i = 0; i < n; i++) 10 { 11 for (j = i; j < n; j++) 12 { 13 if (a[i] < a[j]) 14 { 15 count = count + 1; 16 } 17 } 18 if (count > max) 19 { 20 max = count; 21 } 22 } 23 printf("%d", max); 24 return 0; 25 } </pre>
--	---

Figure 9: Case 3 illustrating the repair of assignment statement of variable and array identifier

Case 4: This case is similar to case 2 but deals with assignment of a variable to another variable instead of array element. In this Figure 10, in the left side of the figure, the lvalue variable "z" inside the For statement is undeclared, and its non-terminal node is ID, it is assigned to the integer type of the variable "i" whose non-terminal is ID.

<pre> 1 int main() { 2 int n, i, j, k; 3 scanf("%d", &n); 4 for(i=1; i<=n; i++) 5 { 6 for(j=1, z=i; j<=i; j++, k--) 7 { 8 if((k%2) == 0) 9 printf("*"); 10 } 11 } 12 return 0; 13 } </pre>	<pre> 1 int main() 2 { 3 int z; 4 int n; 5 int i; 6 int j; 7 int k; 8 scanf("%d", &n); 9 for (i = 1; i <= n; i++) 10 { 11 for (j = 1, z = i; j <= i; j++, k--) 12 { 13 if ((k % 2) == 0) 14 printf("*"); 15 } 16 } 17 return 0; 18 } </pre>
---	---

Figure 10: Case 4 illustrating the fix of variable "z" from the for loop statement

Case 5: This case deals with binary operation involved in an assignment expression statement. In Figure 11, the terminal variable "t" is undeclared and is assigned to the type of the terminal variable "summation" which is of type double. In the statement **summation = summation + t*delx**, there is non-terminal Assignment and its children

nodes being the non-terminal ID with terminal "summation" variable and the node BinaryOp with its corresponding children nodes ID:summation, BinaryOp:+, BinaryOp:* with its children ID:t, ID:delx.

<pre> 1 double sum(double a, double n, double delx) 2 { 3 double summation=0; 4 int j; 5 for (j=0;j<n;j++) 6 {double x=a+j*delx; 7 double r=fabs(f(x)-g(x)); 8 summation=summation+t*delx; 9 } 10 return summation; 11 }</pre>	<pre> 1 double sum(double a, double n, double delx) 2 { 3 double t; 4 double summation = 0; 5 int j; 6 for (j = 0; j < n; j++) 7 { 8 double x = a + (j * delx); 9 double r =fabs(f(x)-g(x)); 10 summation =summation+(t * delx); 11 } 12 13 return summation; 14 }</pre>
---	---

Figure 11: Case 5 demonstrating the error in binary operation and undeclared "t" getting fixed

Case 6: This case is an exact opposite of case 2 where the lvalue in the Figure 12 is an array identifier "b" undeclared, whose non-terminal node is ID and its parent node is ArrayRef and rvalue is a terminal variable "count" whose non-terminal node ID is assigned as integer.

<pre> 1 int main() 2 { 3 int i,j,n,k,count=0,max; 4 scanf("%d",&n); 5 int a[n]; 6 for(i=0;i<n;i++){ 7 scanf("%d",a[i]); 8 } 9 for (i=0;i<n;i++){ 10 for (j=i;j<n;j++){ 11 if (a[j]>a[i]){ 12 count++; 13 } 14 } 15 b[i]=count; 16 count=0; 17 } 18 return 0; 19 }</pre>	<pre> 1 int main() 2 { 3 int b[1000]; 4 int i; 5 int j; 6 int n; 7 int k; 8 int count = 0; 9 int max; 10 scanf("%d", &n); 11 int a[n]; 12 for (i = 0; i < n; i++) 13 { 14 scanf("%d", a[i]); 15 } 16 for (i = 0; i < n; i++) 17 { 18 for (j = i; j < n; j++) 19 { 20 if (a[j] > a[i]) 21 { 22 count++; 23 } 24 } 25 b[i] = count; 26 count = 0; 27 } 28 return 0; 29 }</pre>
---	--

Figure 12: Illustration of case 6 marked by red line indicating the error in assignment statement

Case 7: This case is slightly similar to case 5 but it does not involve any assignment operation. The type can be assigned to a variable not only from lvalue but also from its neighbouring variables involved in a binary operation. As we can see in left side of the Figure 13, the terminal variable "diff" being undeclared is involved in a binary operation BinaryOp:* and BinaryOp:+ with the terminal variables "key" and "a" respectively, so the variable "diff" is assigned to integer type.

```

1 int main()
2 {
3     const double E=0.000001;
4     double a,b,inter , subarea=0;
5     int n,key=0;
6     scanf("%lf%lf%d",&a,&b,&n);
7     inter=(b - a)/n;
8     while(key<n&&diff*key+1< E)
9     {
10         subarea+=1;
11         key++;
12     }
13     return 0;
14 }

```

```

1 int main()
2 {
3     int diff;
4     const double E = 0.000001;
5     double a;
6     double b;
7     double inter;
8     double subarea = 0;
9     int n;
10    int key = 0;
11    scanf("%lf""%lf""%d", &a, &b, &n);
12    inter = (b - a) / n;
13    while((key < n)&&((diff*key)+ 1)<E))
14    {
15        subarea += 1;
16        key++;
17    }
18    return 0;
19 }

```

Figure 13: Demo of case 7 involving the error in while loop statement

Case 8: In this case, the type of a variable is bound from the type of a function call in a conditional expression. In Figure 14, the terminal variable "k" is undeclared in the "for" expression $k \geq \text{hanoi}(j)-1$ and it is being involved in a binary operation $\text{BinaryOp}:\geq$ with the function call $\text{hanoi}(j)$ where the corresponding non-terminal node of the terminal "hanoi" is ID and parent node being FuncCall, is integer.

```

1 int main() {
2     int t,i,n,j;
3     int x;
4     scanf("%d",&t);
5     for(i=1;i<t;i++)
6     {
7         scanf("%d",&n);
8         for(j=0;k>=hanoi(j)-1;j++)
9         {
10             if(hanoi(j)-1==k)
11                 printf("yes");
12             else
13                 printf("no");
14         }
15     }
16     return 0;
17 }

```

```

1 int main()
2 {
3     int k;
4     int t;
5     int i;
6     int n;
7     int j;
8     int x;
9     scanf("%d", &t);
10    for (i = 1; i < t; i++)
11    {
12        scanf("%d", &n);
13        for (j = 0;k>=(hanoi(j) - 1);j++)
14        {
15            if ((hanoi(j) - 1) == k)
16                printf("yes");
17            else
18                printf("no");
19        }
20    }
21    return 0;
22 }

```

Figure 14: Case 8 indicating the undeclared "k" in for loop statement

Case 9: This case is similar to case 8, however instead of a conditional expression with a binary operation, the type of the function call is a rvalue is bound to a lvalue variable in an assignment expression statement with a binary operation. This can be seen from Figure 15, where variable "y" is undeclared in the assignment expression $y = \text{tower}(j)-1$ and is assigned to the type of the function call $\text{tower}(j)$ whose non-terminal node is ID and parent node is FuncCall.

	Identified	Not Identified	Correctly Identified (True Positive)	Wrongly Identified (False Positive)	Correctly Identified + Correct Type Inferred (Fixed)	Wrongly Identified + Wrong Type Inferred (Not Fixed)	Total
Undeclared Variables and Arrays	887(83.7%)	172	857(80.9%)	202	844(79.7%)	215	1059
Undeclared variables - Main function	N/A	N/A	566(99.1%)	5	560(98%)	11	571
Undeclared variables - Multiple functions	N/A	N/A	179(91.7%)	16	172(88.2%)	23	195
Undeclared Arrays - Main functions	N/A	N/A	90(96.8%)	3	90(96.8%)	3	93
Undeclared Arrays - Multiple functions	N/A	N/A	22(78.5%)	6	22(78.5%)	6	28

Table 1: Analysis results of both the undeclared variables and arrays

```

1 int main()
2 {int i,n,j,t;
3  scanf("%d\n",&n);
4  for(i=1;i<=n;i++)
5  {
6    scanf("%d\n",&t);
7    for(j=1;j<=200;j++)
8    {
9      y=tower(j)-1;
10   }
11 }
12 return 0;
13 }

```

```

1 int main()
2 {
3  int y;
4  int i;
5  int n;
6  int j;
7  int t;
8  scanf("%d\n",&n);
9  for(i=1;i<=n;i++)
10 {
11   scanf("%d\n",&t);
12   for(j=1;j<=200;j++)
13   {
14     y = tower(j) - 1;
15   }
16 }
17 return 0;
18 }

```

Figure 15: Case 9 demonstrating the undeclared identifier "y" in assignment statement with a function call

Table 1 shows the results of analysis obtained after performing the compilation of the programs manually where the first row consists of all the programs (1059) that are containing both undeclared variables and arrays in which our approach has located and identified the undeclared variables in 887(83.7%) programs out of total number of 1059 programs. However, our approach has correctly located and identified them in 857(80.9%) programs, but the repair is performed on 844(79.7%) programs by correctly locating as well as inferring and binding their types. The value of the first column "Identified" in the rest of the rows are not applicable (N/A) as the results correspond to programs where undeclared variables are identified and located. Similarly, the second row shows the results for programs with only undeclared variables and consisting of only single main functions (571) out of the identified programs (887). The

Cases	Brief Description
Case 1	Assignment expression statement with a constant on the right-hand side of the expression
Case 2	Assignment expression statement with an array identifier on the right-hand side and an identifier other than array identifier on the left-hand side of the expression
Case 3	Conditional expression statement with a binary operation between the identifier variables
Case 4	Assignment expression statement with an identifier other than array identifier on the right-hand side of the expression
Case 5	Assignment expression statement with a binary operation between the identifier variables
Case 6	Assignment expression statement with an identifier other than arrays on the right-hand side and an array identifier on the left-hand side of the expression
Case 7	Binary operation between identifier variables in a loop expression statement
Case 8	Conditional expression statement with a binary operation between an identifier and a function call expression
Case 9	Assignment expression statement with a binary operation between an identifier and a function call expression

Table 2: Summary of Type Binding Case Description

third row displays the results of programs with undeclared variables and containing two or more functions including main function (195) out of 887 programs. The fourth row demonstrates the results shown by programs only with errors caused due to undeclared arrays and having one and only main function (93) out of the 887 programs. Finally, the last row illustrates the number of programs which contains undeclared variables and having two or more functions along with main function (22) out of those 887 programs. Table 2 shows the summary of various cases along the rows and its brief description message along the columns through which the type binding is performed before the compile-time.

<pre> 1 int main () 2 { 3 int J; 4 int n; 5 int i; 6 int j; 7 int flag = 0; 8 scanf ("%d", &n); 9 int a[51]; 10 for (i = 0; i < n; i++) 11 { 12 scanf ("%d", &a[i]); 13 } 14 for (i = 0; i < n; i++) 15 { 16 for (j = 0; j < n; J++) 17 { 18 if (a[i] == a[j]) 19 { 20 printf ("YES"); 21 flag = 1; 22 break; 23 } 24 } 25 } 26 return 0; 27 }</pre>	<pre> 1 int main () 2 { 3 int l; 4 double a; 5 double b; 6 double k; 7 double p; 8 int n; 9 scanf ("%f" "%f" "%d", &a, &b, &n); 10 k = ((a - b) * 1.0) / n; 11 for (l = 1; l <= n; l++) 12 { 13 if ((l * k) < (-1)) 14 p += k; 15 16 if ((l * k) >= -1) && (l * k <= 1)) 17 p = p + (((l * k) * (l * k)) * k); 18 19 if ((l * k) > 1) 20 p = p + ((l * k) * (l * k)) * (l * k) * k; 21 } 22 printf ("%4f", p); 23 return 0; 24 }</pre>
---	---

Figure 16: Picture on left illustrates repair that caused infinite loop due to variable "J" incremented in loop statement and the right side depicts repair caused by binding type int instead of double

5 Discussion

There are few limitations in our approach. Fixing the undeclared variables that had been due to imperceptible spelling mistakes or a variable that is used only once throughout the program may cause the program to run in an infinite loop or lead to some possible run-time errors. As seen in the left side of Figure 16, instead of incrementing the variable "j" inside the for loop, the programmer had used "J" instead which had caused the program to run into an infinite loop. Another major limitation is in the type binding approach on the right side of Figure 16 shows an example where, the type has been wrongly bound to the variable "l" because "l" is used in the for loop expression in which "l" is assigned to constant "1" as well as it is used in an assignment expression inside the for loop body, in the statement `if((l*k) < -1)`, variable "k" is of the type `double` and "*" is a binary operation, so the variable "j" should be assigned of the type `double` instead it is inferred as an `integer` type due to the former case.

We had seen in our model that a vocabulary in the form of hash table is used for training purposes. The purpose of training neural networks on the hash table is due to the fact that it can be used for recognizing input patterns (keys) in the hash table and can be used to predict the sequences (values). Consider the case where an input pattern is not present in the hash table and we need to predict the sequence, a hash table would have return null in this case but neural networks will give the closest sequence prediction.

The benefits of our approach lies in the fact that our model could be used in real-time as a tool for any C programming environment online or offline editors in locating, reporting and repairing undeclared identifiers for any C programs. Additionally, our model can be used when there are lack of positive bug-free syntactically correct and executing source program reference examples for buggy source programs. Also, our type binding approach will be applicable even for declared variables.

6 Conclusion and Future Work

In this paper, we had seen different cases of one of the most common semantic error: undeclared variables. We had combined AST and LSTM approaches to extract a set of non-terminal and terminal nodes to carry out the classification and prediction tasks of the undeclared variables. We had also seen the generation of clean and buggy-free source programs by performing AST transformation and serialization as well as deserialization of AST to JSON and vice-versa. Furthermore, in this paper we had coined a new term known as Pre-Compile Time Type binding where we had implemented the fix of the types of undeclared variables by binding them their corresponding types before providing it for the compiler to compile them. By our approach, we had correctly identified 81% of the programs that contains only undeclared identifier errors. Also, we had fixed those undeclared identifier errors by binding their corresponding types in 80% of the programs.

In future, we would like to perform automatic repair on different types of syntactic, semantic errors and logical errors. Further, we plan to perform type binding for the limitation cases in Figure 16 as well as also implement a repair approach for the logical errors that arises after the repair of syntactic and semantic errors caused by variables used only once in the program or due to spelling mistakes.

References

- [1] G. Fischer and J. Lusiardi and J. Wolff von Gudenberg, Abstract Syntax Trees - and their Role in Model Driven Software Development. *International Conference on Software Engineering Advances (ICSEA 2007)*. Aug. 2007, pp. 38-38. doi: 10.1109/ICSEA.2007.12.
- [2] Wile, David S, Abstract syntax from concrete syntax. *ICSE*, vol. 97, pp. 472-480, 1997, Citeseer
- [3] Kosovan, Sofia and Lehmann, Jens and Fischer, Asja, Dialogue response generation using neural networks with attention and background knowledge. *Proceedings of the Computer Science Conference for University of Bonn Students (CSCUBS)*, vol. 2017, 2017
- [4] Choudhary, Himanshu and Pathak, Aditya Kumar and Saha, Rajiv Ratan and Kumaraguru, Ponnurangam, Neural Machine Translation for English-Tamil. *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pp. 770-775, 2018
- [5] Bahdanau, Dzmitry and Cho, Kyunghyun and Bengio, Yoshua, Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014
- [6] Luong, Minh-Thang and Pham, Hieu and Manning, Christopher D, Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015

- [7] Goodfellow, Ian and Pouget-Abadie, Jean and Mirza, Mehdi and Xu, Bing and Warde-Farley, David and Ozair, Sherjil and Courville, Aaron and Bengio, Yoshua, Generative adversarial nets. *Advances in neural information processing systems*, pp.2672–2680, 2014
- [8] Rajeswar, Sai and Subramanian, Sandeep and Dutil, Francis and Pal, Christopher and Courville, Aaron, Adversarial generation of natural language. *arXiv preprint arXiv:1705.10929*, 2017
- [9] Guo, Jiaxian and Lu, Sidi and Cai, Han and Zhang, Weinan and Yu, Yong and Wang, Jun, Long text generation via adversarial training with leaked information. *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018
- [10] Liu, Chang and Wang, Xin and Shin, Richard and Gonzalez, Joseph E and Song, Dawn, Neural code completion. 2016
- [11] Li, Liuqing and Feng, He and Zhuang, Wenjie and Meng, Na and Ryder, Barbara, Cclearner: A deep learning-based clone detection approach. *2017 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pp. 249–260, 2017, IEEE
- [12] Dam, Hoa Khanh and Tran, Truyen and Pham, Trang and Ng, Shien Wee and Grundy, John and Ghose, Aditya, Automatic feature learning for vulnerability prediction. *arXiv preprint arXiv:1708.02368*, 2017
- [13] Harer, Jacob and Ozdemir, Onur and Lazovich, Tomo and Reale, Christopher and Russell, Rebecca and Kim, Louis and others, Learning to repair software vulnerabilities with generative adversarial networks. *Advances in Neural Information Processing Systems*, pp. 7933–7943, 2018
- [14] Ahmed, Umair Z and Kumar, Pawan and Karkare, Amey and Kar, Purushottam and Gulwani, Sumit, Compilation error repair: for the student programs, from the student programs. *2018 IEEE/ACM 40th International Conference on Software Engineering: Software Engineering Education and Training (ICSE-SEET)*, pp. 78–87, 2018, IEEE
- [15] Gupta, Rahul and Pal, Soham and Kanade, Aditya and Shevade, Shirish, Deepfix: Fixing common c language errors by deep learning. *Thirty-First AAAI Conference on Artificial Intelligence*, 2017
- [16] Bhatia, Sahil and Singh, Rishabh, Automated correction for syntax errors in programming assignments using recurrent neural networks. *arXiv preprint arXiv:1603.06129*, 2016
- [17] Jayanthi, R and Florence, Lilly, Software defect prediction techniques using metrics based on neural network classifier. *Cluster Computing*, pp. 1–12, 2018, Springer