# Code Summarization using Machine learning

Musthaq Ahamad, P Harikrishnan, Lalith Sagar Manoj, and Maneesh VS

Computer Science and Engineering Dept., Sahyadri College of Engineering & Management, Adyar, Mangaluru-575007

**ABSTRACT**

**Programming language processing (like normal language processing) is a hot research theme in the world of software engineering community; it has additional exciting developing interest in the AI community or society. Be that as it may, not quite the same as a natural language sentence, a program contains rich, express, and muddled basic data. Henceforth, standard NLP models may not be feasible for projects. In this paper, we propose a novel tree-based convolutional neural network (TBCNN) for programming language handling, in which a convolution piece is structured over project's abstract syntax trees to catch basic data. TBCNN is a conventional engineering for programming language processing; our tests demonstrate its viability in two diverse program examination according to functionality: arranging programs as indicated by programming functionality, and distinguishing code snippet of specific pattern. TBCNN beats standard techniques, counting a several neural models for NLP.**

*Keywords:*
NLP, TBCNN, PROGRAMMING LANGUAUE PROCESSING, AST,
SOFTWARE ENGINEERING.

## 1.  INTRODUCTION

Analysts from different society in computer science takes special interest in applying Artificial Intelligence consciousness (AI) systems to software Engineering (SE) building issues(Dietz et al. 2009; Bettenburg and Begel 2013; Hao et al. 2013). In the study of SE, dissecting program source code, called programming language summarizing, in this paper show specific significance. Despite the fact that computers can run programs which they are fed, they don't genuinely understand the programs. Dissecting source code gives a method for assessing projects' conduct, usefulness, multifaceted nature, and so forth. For example, consequently distinguishing source code bits of specific examples help developers to find surrey or wasteful calculations in order to improve code quality. Another model is overseeing expansive programming vaults, where programmed source code grouping and labeling are essential to programming reuse. Programming language handling, indeed, fills in as an establishment for some, SE assignments, e.g., requirement analysis (Ghabi and Egyed 2012), programming improvement and support (Bettenburg and Begel 2013).

Hindle et al. (2012) exhibit that programming dialects, like natural languaues, additionally contain copious measurable properties, which are significant for program examination. These properties are hard to catch by people, in any case, legitimize learning-based methodologies for programming language handling. Be that as it may, existing ML program processing depends to a great extent on feature building, Which is work concentrated and specially appointed to a particular undertaking when needed, e.g., code clone detection (Chilowicz, Duris, and Roussel 2009), and bug recognition (Steidl and Gode 2013). Further, proof in the ML writing proposes that human engineered funtionalities may neglect to catch the idea of information or its nature, so they might be far more atrocious than consequently learned ones.

The profound neural system, otherwise called deep learning, is an exceedingly mechanized learning machine. By investigating different layers of non-direct automatic learning machine, the profound design can consequently learn confused hidden features in the data, which are urgent to the assignment of intrigue. Over the recent years, deep learning has made huge leaps forward in different fields, for example, discourse acknowledgment (Dahl, Mohamed), computer vision (Krizhevsky, Sutskever and silvia 2006), and natural language processing (Collobert and Weltom 2007).

In spite of certain common factors between natural language dialects and programming languaue dialects, there are additionally clear contrasts (Sheet, Ratanamahatana, and Myers 2001). In light of a formal language, programs contain rich and unambiogous basic data. Despite the fact that structures additionally exist in natural language, they are not as stringent as in projects. linker (1993) represents an intriguing precedent, "The dog the stick the flame consumed beat bit the kitten." This sentence consents with all syntax rules, yet such a

large number of attributive clauses are settled. Subsequently, it can barely be comprehended by individuals because of the constraint of human instinct limitation. On the opposite, three settled circles are normal in projects. The parse tree of a program, indeed, is commonly a lot bigger than that of a natural languaue sentence—there are roughly 190 nodes by and large in our trial, while a sentence contains just 20 words in an assessment investigation dataset (voncher 2013). Further, the language structure rules "alias" neighboring relationship among program parts. The announcements inside and outside a circle, for instance, try not to shape one semantic gathering, and in this way are not semantically neighboring. On the above premise, we think more powerful neural models are in need to catch auxiliary data in projects.

In this paper, we propose a novel Tree-Based Convolutional Neural Network (TBCNN) in view of projects' Abstract structure trees (ASTs). We additionally present the idea of "continuous binary tree" and apply dynamic pooling to adapt to ASTs of various sizes and shapes. The TBCNN model is a conventional design, and is connected to two SE assignments in our trials—characterizing programs by functionalities also, recognizing code pieces of specific examples. It outflanks standard techniques in the two errands, including the recursive neural system (voncher 2013) proposed for NLP. As far as we could possibly know, this paper is likewise the first to apply deep neural systems to the field of programming language processing.

## 2. RELATED WORK

### 2.1 Research and experiments

Deep neural systems have made noteworthy achievements in numerous fields. Stacked confined Boltzmann machines and autoencoders are effective pretraining techniques (Riyadh Mehrez, Osindero, and Teh 2006; Bengio et al. 2006). They investigate the hidden highlights of information in an unsupervised way, furthermore, give an increasingly significant introduction of loads for some other time regulated learning with profound neural systems. These methodologies function admirably with conventional information (for example information situated in a complex implanted in a specific dimensional space), yet they may not be appropriate for preparing programming language, since projects contain rich auxiliary data. Further, AST structures likewise shift to a great extent among various information tests (programs), and consequently they can't be sustained legitimately to a fixed-estimate arrange.

To catch unambigous structures in information, it might be significant furthermore, useful to incorporate human priors to the systems (Bernardo Cavalhio Silva, Kyle Walker, and Vincent Kompany 2013). One model is convolutional neural systems (CNNs, LeCun et al. 1995; Krizhevsky, Heu Min Son, and Hinton 2012), which determine spatial neighboring data in information. CNNs work with signs of a specific measurement (e.g., pictures); they additionally come up short to catch tree-basic data as in projects. voncher *2013* (2013, 2011b) propose a recursive neural system (RNN) for NLP. Albeit auxiliary data can be coded somewhat in RNNs but the real disadvantage is that just the root highlights are utilized for directed learning, which covers lighting up data under a muddled neural engineering. RNNs likewise experience the ill effects of the trouble in preparing because of the long reliance way amid back-engendering (dengio, Simare, and Frasloni 1992)

### 2.2 Subsequent work

After the primer variant of this paper was preprinted on arXiv, 2 Zaremba and Sutskever (2014) utilize intermittent neural systems to evaluate the yield of limited python programs. Piech et al. (2015) form recursive arranges on Hoare triples. With respect to proposed TBCNN, we extend it to process syntactic parse trees of normal dialects (Mou et al. 2015); Duvenaud et al. (2015) apply a comparable convolutional organize over charts to break down particles

### 2.3 .Tree-based Convolutional Neutral Network

Programming dialects have a characteristic tree portrayal dynamic language structure tree (AST). Figure 1a demonstrates the AST of the code bit "int a=b+3;".3 each hub in the AST is a dynamic segment in program source code. A hub p with youngsters speaks the developing procedure of the segment p→c1…….cN demonstrates the general engineering of TBCNN. In our model, an AST hub is first spoken as an appropriated, genuine esteemed vector so that the (unknown) highlights of the images are caught. The vector portrayals are found out by a coding rule in our past work (Peng et al. 2015). At that point we structured a lot of subtree including indicators, called the tree-based convolution portion, sliding over the whole AST to remove basic data of a program. We from that point apply dynamic pooling to accumulate data over distinctive pieces of the tree. At last, a shrouded layer and a yield layer are included. For regulated grouping undertakings, the actuation capacity of the yield layer is softmax. From this segment, we initially clarify the coding paradigm for AST hub, picking up portrayal, filling in as a pretraining period to prepare programming language. We at this point describes the proposed TBCNN model, including a coding layer, a convolutional layer, and a pooling layer. We moreover use extra data on managing hubs that have changing quantities of tyke hubs, as in ASTs, by presenting the concept of ceaseless paired trees organized over diagrams to break down atoms

### 2.4 Representation Learning for AST Nodes

Vector portrayals which is known as embeddings, can catch fundamental implications of discrete images, as AST hubs. We propose in our past work (Peng et al. 2015) an unsupervised way to deal with learn program vector portrayals by a coding paradigm, which is filled in a way of pretraining.
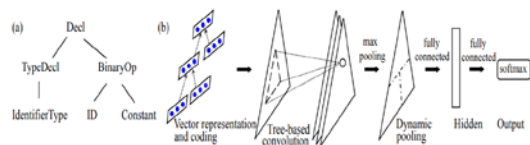
Figure 1: (a) Illustration of an AST, corresponding to the C code snippet "int a=b+3;" It should be notice that our model takes as input the entire AST of a program, which is typically much larger. (b) The architecture of the Tree-Based Convolutional Neural Network (TBCNN). The main components in our model include vector representation and coding, tree-based convolution and dynamic pooling; then a fully-connected hidden layer and an output layer (softmax) are added.

A conventional model for portrayal learning is "smoothness"— in wich images have comparative component vectors (Virgil Van Dijik, and Lucas Moura 2013). For instance, the images are comparable on the grounds that are identified with control stream, especially circles. However, they are unique in relation to ID, since ID likely speaks few information. In our case, we might want the youngster hubs' portrayals to "code" their parent hub's by means of a solitary neural layer, along with both vector portrayals and coding loads are found out. Formally, let vec(.) belongs to RNf be the component portrayal of an image, where $N_f$ is the component measurement. For each non-leaf hub p and its immediate children c1… cn, we want your manuscript in different sections.

### 2.5 Dynamic Pooling

Removed, what's more, another tree is created. The new tree has precisely a similar shape and size as the first one, which is changing among various projects. So, the separated highlights can't be used as input to a fixed-estimate neural layer. Dynamic pooling (van dijik et al. 2011a) is used to tackle this issue. The easiest methodology, maybe is to pool all highlights to one vector. We call this single direction pooling. Solidly the most upper value in each measurement is taken from the highlights that are identified by tree-based convolution. We likewise propose an option, three-way pooling, where highlights are pooled to 3 sections, TOP, LOWER LEFT, and LOWER RIGHT, concurring to the their location in the AST (Figure 2b). As we will see from the test results, the straightforward single direction pooling works the same way a three-way pooling works. Thus we used single direction pooling in our observations. Post pooling, the highlights are completely associated with a concealed layer and afterward given as input to the retuen layer for. With the dynamic pooling process, auxiliary highlights along the whole AST achieve the yield layer with short ways. Henceforth, they can be prepared successfully by back-spread.

### 2.6 The "Continuous Binary Tree" Model

As expressed, one issue of coding and convolving is that we can't decide the quantity of weight grids because AST hubs have distinctive quantities of child nodes. One conceivable arrangement is the successive group of-words model (CboW,

Vijoven , 2012),5 however position data will be lost fully. Such methodology is likewise utilized in Harmann and Blunson (2012). Boncer . (2013) distribute an alternate weight lattice as parameters for each location; yet, this technique neglects to grow since there will be an immense number of various positions in ASTs

In our model, we see any subtree as a "binary" tree, without considering its size and shape. That is, we have just three weight grids as parameters for convolution, and two for coding. We consider it a successive tree.

## 2 EXPERIMENTS

We initially survey the vector portaryals both subjectively and quantitatively. After this we asses TBCNN in two regulated learning assignments, and lead model investigation. The dataset of out trails originates from an educational programming open judge (OJ) system. There are a huge number of programming issues on the OJ framework understudies present their source code as the answer for a certain issue; The OJ framework naturally makes a decision about the legitimacy of submitted source code by running the program. We downloaded the source code and the relating programming issues (spoke to as IDS) as our dataset

### UNSUPERVISED PROGRAM VECTOR REPRESENTATIONS

We applied the codding criterion of pretrainning to all C code in the OJ system, and getting as result AST nodes vector representations.

**Qualitative analysis**. Figure 3a outlines the progressive Grouping result dependent on a subset of AST hubs. As illustrated, the images fundamentally fall into three classes: (1) BinaryOp, ArrayRef, ID, Constant are gathered since they are identified with information reference/control; (2) For, If, While are comparable since they are identified with control stream; (3) ArrayDecl, FuncDecl, PtrDecl are comparable since they are statements. The outcome is very reasonable since it is steady with human comprehension of programs.

**Quantitative analysis:** we likewise assessed pretraining's Impact on regulated learning by encouraging the assessed presentation to a program characterization task. Figure 3b plots the expectations to absorb information of both preparing and approval, which are contrasted and arbitrary instatement. Unsupervised vector portrayal learning quickens the administered preparing process by almost 1=3, demonstrating that pretraining captures basis highlights of AST hubs, and that they can rise abnormal state includes unexpectedly amid regulated learning. In any case, pretraining has a restricted

impact on the last exactness. One conceivable clarification is that the quantity of AST hubs is little: the pycparser, we use, recognizes just 44 images. Subsequently, their portrayals can be enough tuned in a directed manner.
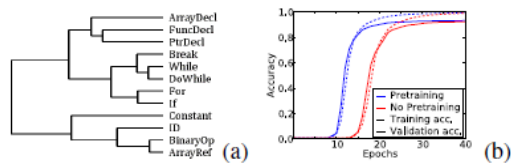


Figure 3: Analysis of vector representations. (a) Hierarchical clustering based on AST nodes' vector representations. (b) Learning curves with and without pretraining.

| Statistics | Mean | Sample std. |
|---|---|---|
| # of code lines | 36.3 | 19.0 |
| # of AST nodes | 189.6 | 106.0 |
| Average leaf nodes' depth in an AST | 7.6 | 1.7 |
| Max depth of an AST | 12.3 | 3.2 |

Table 1: Statistics of our dataset.

| Hyperparameter | Value | How is the value chosen? |
|---|---|---|
| Initial learning rate | 0.3 | By validation |
| Learning rate decay | None | Empirically |
| Embedding dimension | 30 | Empirically |
| Convolutional layers' dim. | 600 | By validation |
| Penultimate layer's dim. | 600 | Same as conv layers |
| $l_2$ penalty | None | Empirically |

Table 2: TBCNN's hyperparameters.

We think the pretraining criterion is useful and advantages for TBCNN, because teaching deep neural networks is a long process, specifically when aping hyperparameters. The rethought vector representations are used continuously in our experiments below.

**Classifying Programs by Functionalities**
**Task description**

In software engineering, grouping programs by the function it performs is an essential issue for diffrent software advancement tasks. For ex, in a huge software hold (e.g., SourceForgge), software items are mostly arranged into brackets, a usual specification for which is by tasks done. With program specification, it becomes easier to instinctively tag a software part newly entered into the repository, which is advantageous for software re use during the growth process. In our experiment, we used TBCNN to group source code in the OJ system. The target identification of a data sample is one of 104 programming issues (shown as an ID). That is, programs with a similar target identification have the same use. We randomly select accurately 500 programs in each section, and thus 52,000 samples in all, which were again randomly split by 3:1:1 for training, validation, and testing. Statistics related are shown in Table 1**.**

**Hyperparameters** TBCNN's hyperparameters are given

In Table 2. The methods we use contain SVM and a deep feed-forward neural network based on hand-crafted features, namely bag-of-words (BoW, the counting of each symbol) or bag-of-tree (BoT, the counting of 2-layer subtrees). We also asses our model with the recursive neural network

(RNN, Socher et al. 2011b). Hyperparameters for baselines are listed as follows.

**SVM**. The linear SVM has one hyperparameter C; RBFSVM has two, C and. They are tuned by validation overthe set { ; 1; 0:3; 0:1; 0:03; ….. } with grid search.

**DNN**. We applied a 4-layer DNN (including input) empirically. The hidden layers' dimension is 300, chosen fromf100; 300; 1000g; learning rates are 0.003 for BoW and 0.03for BoT, chosen from {:003; ……; 0:3} with granularity 3x.`2 regularization coefficient is $10^{-6}$ for both BoW and BoT,chosen from$\{10^{-7}, ………, 10^{-4}\}$ with granularity 10x, and also no regularization.

**RNN**. Recursive units are 600-dimensional, as in
Our method. The learning rate is chosen from the set
{…. 1:0; 0:3; 0:1……}, and 0.3 yields the highest validation performance.

# 3 CONCLUSION

In the given study we clustered vivid neural systems to the field of programming language processing. Because of the rich and unambiguous tree structures of this projects, we proposed the novel Tree-Based Convolutional Neural Network (TBCNN). In our proposed model, program vector portrayals are found out by the coding basis; basic highlights are identified by the convolutional layer; the consistent double tree and dynamic pooling authorize our model to adjust to trees of changing sizes and shapes. Trial results represents the widespread our model to pattern strategies.

# 4 ACKNOWLEDGMENT

# 5 REFERENCES

[1] Hector Bellerin, Y.; Dembele, P.; Mauricio Pochettino, D.; and Jesus Perez, H.2006. Greedy layer-wise training of deep networks. In NIPS.

[2] Hugo Lloris, Y.; Kieran Trippier, A.; and Victor Vanayama, P. 2013. Representationlearning: A review and new perspectives. IEEE Trans.Pattern Anal. Mach. Intell. 35(8):1798–1828.

[3] Fernando Llorente, Y.; Dele Ali, P.; and Frasconi, P. 1994. Learninglong-term dependencies with gradient descent is difficult. IEEETrans. Neural Networks 5(2):157–166.

[4] Paulo Gazzaniga, N., and Christian Erickson, A. 2013. Deciphering the story of software development through frequent pattern mining.In ICSE, 1197–1200.

[5] Dany Rose, M.; Duris, E.; and Jan Vertonghen, G. 2009. Syntax tree fingerprinting for source code similarity detection. In Proc.IEEE Int. Conf. Program Comprehension, 243–247.

[6] kWEV, R., and Weston, J. 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning. In ICML.

[7] Kevin Debruyna, R.; Weston, J.; Ederson, L.; Laporte, M.; Otamendi, K.; and Kuksa, P. 2011. Natural language processing (almost) from scratch. JRML 12:2493–2537.

[8] Dahl, G.; Fernandhino, A.; and Hinton, G. 2010. Phone Recognition with the mean-covariance restricted Boltzmann machine. In NIPS.

[9] Dietz, L.; Dallmeier, V.; Zeller, A.; and Gundagoan, T. 2009.Localizing bugs in program executions with graphical models. In NIPS.

[10] Phil Foden, D.; Maclaurin, D.; Sergio-Kun-Aguero, J.;Andre Gomez, R.; David Silva, T.; Rahim Sterling, A.; and Adams, R. 2015. Convolutional networks on graphs for learningmolecular fingerprints. arXiv preprint arXiv:1509.09292.

[11] Gabrial Jesus, A., and Egyed, A. 2012. Code patterns for automatically validating requirements-to-code traces. In ASE, 200–209.

[12] David De Gea , D.; Lan, T.; Romero, H.; Guo, C.; and Lindolf, L. 2013.Is this a bug or an obsolete test? In Proc. ECOOP, 602–628.

[13] Smalling, K., and Luke Shaw, P. 2014. Multilingual models for compositional distributed semantics. In ACL, 58-68.

[14] Ashley Young, A.; Barr, E.; Su, Z.; Valencia, M.; and Devanbu, P. 2012. On the naturalness of software. In ICSE, 837–847.

[15] Matic, G.; Osindero, S.; and Teh, Y. 2006. A fast learning algorithm for deep belief nets. Neural Computation 18(7):1527–1554.

[16] Paul Pogba, N.; Harrera, E.; and Blunsom, P. 2014. A convolutional neural network for modelling sentences. In ACL,655–665.

[17] Juan Mata, A.; Bernd Leno, I.; and Sokartis, G. 2012. ImageNet classification with deep convolutional neural networks. In NIPS.

[18] Koscieleny, Y.; Ncho Monreal, L.; Mustafi, L.; Brunot, A.; Cortes, C.; Elneny, J.; Drucker, H.; Guyon, I.; Muller, U.; and Mkhitaryan, E.1995. Comparison of learning algorithms for handwritten digit recognition. In Proc. Int. Conf. Artificial Neural Networks.

[19] Eden Hazard, T.; Aarom Ramsey, I.; Chen, K.; Mesut Ozil, G.; and Dean, J. 2013. Distributed representations of words and phrases and their compositionality. In NIPS.

[20] Mou, L.; Xhaka, H.; Li, G.; Xu, Y.; Geundeouzi, L.; and Jin, Z. 2015. Discriminating neural sentence modeling by tree-based convolution. In EMNLP, 2315–2325.

[21] Pane, J.; Lucas Toriera, C.; and Riyadh Mehrez, B. 2001. Studying the language and structure in non-programmers' solutions to programming problems. Int. J. Human-Computer Studies 54(2):237–264.for paraphrase detection. In NIPS.

[22] Alexander Lcazette R.; Aubameyang J.; Huang, E.; Ng, A.; and Manning,C. 2011b. Semi-supervised recursive autoencoders for predicting Sentiment distributions. In EMNLP, 151–161.[2013] Socher, R.; Alex Evobi, A.;Wu, J.; Peter Chech, J.; Manning, C.;Ng, A.; and Potts, C. 2013. Recursive deep models for semantic Compositionality over a sentiment treebank. In EMNLP, 1631–1642.