

---

# Analyze Programming Languages

---

**Yiqiao Yin**  
W.Y.N. Associates, LLC

## Abstract

1        This assignment reports the research findings for analysis of pro-  
2        gramming languages. The goal is to understand basic programming  
3        languages such as C++. There are two parts for this assignment report.  
4        First, the report focuses on understanding the basics of programming  
5        languages. Second, the report summarizes some current research built  
6        using C++ where the C++ forms the basis of the research (all related  
7        papers are two years recent and from peer-reviewed sources).

## 8    **1    Understanding C++**

### 9    **1.1    Introduction**

10    The modern day advancement in fields such as healthcare, technology, finance, and  
11    even agriculture all require the usage of many computers at different perspectives of  
12    the business parts. It is extremely challenging to operate or manage business without  
13    the participation of computers. Hence, the development of problem-solving skills with  
14    coding is extremely important.

15    Fundamentally, the task of coding by itself is a form of problem-solving. In practice, it  
16    is often time the case to be able to write many different approaches of solutions tackling  
17    one problem. Hence, a flexible coding style is also very important. The understanding  
18    of the programming languages can largely increase a person's logical thinking as well as  
19    execution in a neat and elegant manner.

20    With such motivation discussed as the above, it is important to, first of all, address two  
21    basic types of languages: assembly language and machine language. Assembly language  
22    or ASM is a low-level language where the relationship between the language instruction  
23    and the skeleton of code presentation is strong Juneja (2009); Sinkov (1963). We use  
24    ASM as a replacement for assembly language for the rest of the paper. The style ASM is  
25    short and elegant. It usually consists of a simple statement or execution of code. The  
26    coding style is also short and tight with on an one-to-one basis. For the execution of  
27    ASM, the code is converted into a small program that is called assembler, which matches  
28    the concept where the name ASM suggests. The ASM uses the term assembler to refer  
29    to this unity program that is created by this executable machine. A fun fact is that the  
30    work Wilkes et al. (1951) introduces the term "assembler" in the first time in history.  
31    This concept also lays the ground of today's terminology "assemble machine". It is a  
32    necessary building block to bridge the gap between software programs and the hardware  
33    platforms, because ASM translates high-level language into machine language. The most

34 basic example in a computer for ASM is the binary code consists of zeros and ones. It is  
 35 rare for an organization to put humans directly in charge of writing code such as ASM.  
 36 This is where machine language shows its strength. Machine language is a low-level  
 37 programming language that consists of language instructions Juneja (2009); Sinkov  
 38 (1963). It is used to control the central processing unit or CPU. The instructions to  
 39 control CPU often carries a very particular task. Ideally, the code is written to execute the  
 40 jobs or tasks sent to the CPU as efficient as possible which also requires the knowledge  
 41 of data structures. As opposed to the ASM, machine language can be done by human  
 42 but the management style requires detailed investigation. Machine language allows us to  
 43 build more complicated computer program to send command to CPU inside a computer.  
 44 A comparison is presented between the ASM and machine language in Table 1.

Assembly Language	Machine Language
Between high- and low-level	Low-level language
Uses English alphabet	Uses zeros and ones
Understood by human	Understood by CPU
Memory exists for instructions	No memory

Table 1: **Comparison between ASM and Machine Language.**

## 45 1.2 Famous Applications

46 A famous coding language is C++. In real world application, there are many areas where  
 47 C++ can provide crucial developmental role in the process.

48 The first example is the gaming industry. The design of a video tracks down to the core  
 49 of the hardware, which is something C++ is really good at. The programming language  
 50 is designed to easily operate. The procedure of game programming over CPU consists  
 51 of intensive coding and fast executed coding structure. It is sometimes organized over  
 52 multiple network and can involve 3D design with multi-player setup. Famous games  
 53 such as Witcher 3, Counter-Strike, Warcraft, and Doom are great application products  
 54 using C++.

55 The next field of application is Graphical User Interface or GUI, which is a type of user  
 56 interface that bridge the gaps between user and graphical icons instead of code. It is a  
 57 system of interactive visual components for computer software of which C++ shows  
 58 significant strength at. Famous applications in GUI are Adobe and Windows Media  
 59 Player.

60 For data scientists and senior data scientists, the workload cannot leave the realm of data  
 61 structures. This is an area where database management can be extremely important for  
 62 the role of a data scientist. The C++, in this case, serves as a great tool to develop writing  
 63 of the database management software. Famous examples are MySQL and Postgre.

64 A whole list can be found in Section 1.5 of the work by Juneja (2009). We present a few  
 65 selected applications here.

- 66 1. Computer operational systems such as Windows XP and OS.
- 67 2. Real-time control applications.
- 68 3. 2D and 3D animations and gaming design.
- 69 4. Speech recognition and dictation systems.

### 70 1.3 Execution Diagram

71 The purpose of a compiler is serves as a processor in the creation of the binary code.  
72 The process the procedure of developing a program using a language such as C++ is  
73 summarized below in Figure 1.

74 To understand the executing order of the diagram from the work by Juneja (2009). The  
75 procedure starts with a text editor. This is the place where the code or the English  
76 syntax is typed into a script. This syntax is tailored to each computer language. It is  
77 recommended to have a round of bug check. This is indicated using a diamond shape  
78 box in the diagram. If there is no error, the procedure continues. If there is an error, it  
79 is recommended to go the beginning which is the text editor to rewrite the code which  
80 is the English syntax. Next, a preprocessor is used to go through the source code for  
81 the compiler. Then a compiler takes over to create a binary file for source code. This  
82 source code needs to be linked, which gains its name “linker”. The linker rends the code  
83 and reports bugs when the code has errors present. If that is the case, it is recommended  
84 to go back to the syntax part which is the first step to start over. However, if there is  
85 no bugs reported, the program successfully renders and compiles the code. Then the  
86 program executes the binary file and it is ready to run. The user can then execute the  
87 program and runs the program. There is another round of bug report. If there is any bug  
88 present when running the program, it is recommended to go back to the syntax or text  
89 editor for a third time. If the program is free of bugs, the program runs, executes, and  
90 the results are generated.

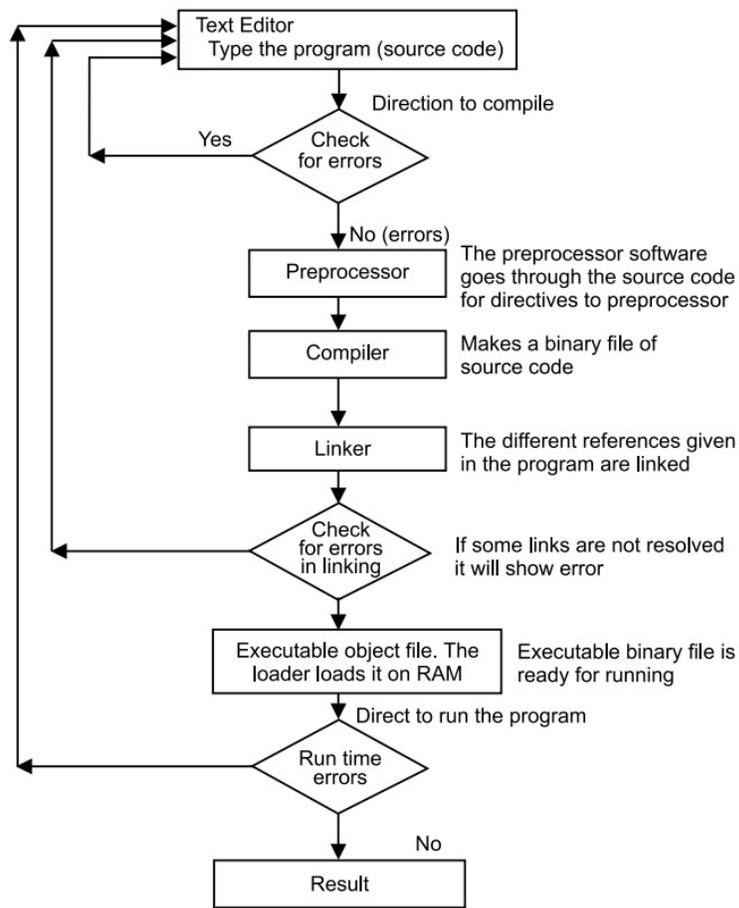
## 91 2 C++ in Research

### 92 2.1 Background

93 This portion of the assignment investigates and reviews the work by Abuhamad et al.  
94 (2021). The field of authorship identification is fundamentally a natural language  
95 processing task. Many former scholar has conducted thorough research in this field Juola  
96 (2008); Kešelj et al. (2003); Koppel et al. (2009); Stamatatos (2009). However, there  
97 is little work done in authorship identification in source code and proposed computer  
98 programs Burrows et al. (2014). Based on the unique stylometric features of an author  
99 or a programmer, source code can be structurally written in a fundamentally different  
100 approach. It is this distinct nature that allows each programmer to tailor their experience  
101 according to the pipeline of the particular coding project. From linguistic perspective,  
102 the notion of stylometric serves a key purpose in this context, because stylometric is  
103 the analysis of features of the writing style that can be statistically and quantitatively  
104 measured which can be a range of features such as sentence length, vocabulary diversity,  
105 frequency, order of phrases, and so on.

106 Their work investigated different stylometry of the writings. This means that experiments  
107 are conducted to compare for the intrinsic writing style in machine learning algorithms.  
108 In this case, the concept of stylometry of stylometric feature is referring to the writing  
109 style of machine learning algorithms or code in general that can be quantified statistically.  
110 These features can refer to sentence length, vocabulary diversity, and frequencies (of  
111 words, word forms, and so on).

112 The motivation for the work by Abuhamad et al. (2021) is to develop a system for  
113 a numerous of applications that can take advance or be built upon the understanding  
114 of authorship syntax and writing style. Software identification is a major component  
115 of software forensics and security analysis. A major benefit is to establish potential



**Fig. 1.4:** Sequence of processes in development of a program in C/C++ language.  
(Also See Appendix E for more details on starting programs in C++)

**Figure 1: Diagram of the Developmental Procedure.**

116 malicious code or even malware that intend to deliberately deliver malicious activity.  
117 Another application can be the adaption and potential development of sophisticated  
118 plagiarism detection Burrows et al. (2007), authorship disputes Wilcox (1998), copyright  
119 infringement, and software integrity investigations Malin et al. (2008).

## 120 **2.2 Review of Literature and Similar Work**

121 To tackle the above problems, recent work in software authorship identification has  
122 proposed several interesting techniques Koppel et al. (2009); Spafford and Weeber  
123 (1993); Burrows et al. (2007); Caliskan-Islam et al. (2015). However, their work cannot  
124 be generalized to other field and it is challenging to apply the same feature extraction  
125 technique to other programming languages such as Java or Python. In other words, the  
126 model built for identifying authorship in C++ should only stay in documents and source  
127 code that is C++ based. A former study proposed writeprints which is a technique that  
128 captures the author's stylometry in diverse corpus Abbasi and Chen (2008). The work by  
129 Uzuner and Katz (2005) conducted comparative study to investigate different stylometric  
130 information in the author's writing to identify the author's attribution and identification.  
131 The work is taken further into the possible attack of cybercriminals by using textual  
132 entries Afroz et al. (2014).

## 133 **2.3 Role of C++ in This Research**

134 The research uses Term Frequency-Inverse Document Frequency or TF-IDF method to  
135 extract features from the code syntax and input the extracted features into a Recurrent  
136 Neural Network (RNN). The authors showed state-of-the-art prediction performance on  
137 a number of experiments. Amongst all 5 listed contributions in the paper, 4 of the data  
138 sets consists of C++ and the language can also be compared with C, Java, and Python.

## 139 **2.4 Research Question**

140 Though the authorship style is unique, the source code is practiced in the following way.  
141 In order to process the source code, a program is compiled to the original source code  
142 and this is in executable binary bits. This is then decompiled to generate pseudo-code.  
143 The syntax at a low-level construction is largely dependent of the language itself and  
144 can impose an inherent inflexibility in the architecture of the code. This nature poses  
145 great challenge for machine learning scientists to learn about the stylometric of the text  
146 using source code. In plain English, there are only so many ways to write a chunk of  
147 code and at a basic level of the code it is common that most people write the code in a  
148 similar way in terms of the style features used in the stylometric analytical system.

149 Several questions raised by the author Abuhamad et al. (2021) that deserve attention are  
150 listed in the following. Is it possible to use deep learning technology to help investigate  
151 or potentially classify the software authors? If yes, how can this experimental procedure  
152 be done? In addition, how can we generalize deep learning based approach to wider  
153 syntax? If deep learning can be used in one type of syntax, can this methodology be  
154 used in other syntax without the prior knowledge? In the case of using this technology in  
155 other languages, can the result still perform robustly or will the performance be damaged  
156 due to new language or syntax environment?

## 157 **2.5 Research Findings**

158 The authors developed a sequential model with Recurrent Neural Network or RNN. The  
159 model aims to extract the features from the syntax database and use these extracted  
160 features to construct the machine learner. For each instance, the authors fed in either  
161 complete or incomplete source code or syntax to create explanatory features. These  
162 features are generally believed to have high quality and contains distinctive authorship  
163 stylometric. Unlike former research that relies on prior knowledge for data transforma-  
164 tion, the proposed work extract features automatically using RNN-related deep learning  
165 methods. The author proposes to use Term Frequency-Inverse Document Frequency  
166 (TF-IDF) Burrows and Tahaghoghi (2007); Frantzeskou et al. (2006); Kešelj et al. (2003).  
167 Since TF-IDF is non-bayesian, it does not require any prior knowledge of any particular  
168 programming language and nor does it need any high-level translation. Hence, the author  
169 believed that it is a more resilient approach to language specifics. It is also beneficial  
170 because it can be generalized language-wise when the source code is accessible at hand  
171 and the program can be processed into binary format. The important features that help  
172 explain the predictions, after the model is trained and fitted on the training dataset, are  
173 famous keywords of the used programming language. This result poses an interesting  
174 complexity, because this indicates that the authorship cannot be easily identified by  
175 variable names but by their dramatically different programming style. The authors are  
176 able to achieve the state-of-the-art accuracy using the extracted features.

177 The authors provided a sequential model that is RNN based and delivered the state-  
178 of-the-art prediction performance using features extracted from TF-IDF. Over large  
179 dataset of 150 C++ programmers, the prediction performance is near perfect and at  
180 a 99% accuracy. An even larger dataset with a huge number of programmers (8,903  
181 programmers) is tested using the same proposed approach and the accuracy is 92.3%.

182 In addition, to ensure the features extracted are meaningful, the authors also used a  
183 random Forest Classifier or RFC as a classifier to learn from the TF-IDF-based features.  
184 The RFC reads in the features and then build a classifier by splitting different trees. The  
185 number of trees is a tuning parameter. Each tree behaves like a decision tree and grows  
186 many branches. Each branch is a split from comparing the feature at that branch with a  
187 particular threshold. An advantage of the proposed model is that it requires only seven  
188 files for each programmer in a dataset that has 5.5 times the number of programmers  
189 than that of the prior work and that it is able to create similar performance.

190 The third contribution from Juneja (2009) is that analysis is tested on different pro-  
191 gramming languages such as C++, C, Java, and Python. Analysis is also done using  
192 combination of different programming languages (C++/C, C++/Java, and C++/Python).  
193 The purpose is to test out whether the proposed model can be generalized enough for  
194 multiple different computer languages.

## 195 **2.6 Limitations and Future Research**

196 There are some limitations to the author's paper and these limitations lead to future  
197 research opportunities. This report summarizes below.

198 First, the experiments are conducted under the assumption that each source code sample  
199 is finished with a single programmer. In other words, in a tabular data form, each row can  
200 be considered as a single observation or an instance where this observation happens to  
201 be only a single programmer. There assumption dictates that there is no situation where  
202 more than one programmers collaborate and finish the source code sample together.  
203 This assumption rarely holds true in practice, because the common practice is a team

204 of programmers collaborate and check each others' code. In that case, this assumption  
205 does not hold true and we need to release the assumption by accepting the situation that  
206 collaboration can happen in the industry.

207 Second of all, the experiments assumes that the authorship and stylometric of every  
208 programmer is unique and poses great differences amongst each other. In addition, the  
209 experiment assumes that the programming styles do not vary sample by sample on a  
210 single programmer. These assumptions rarely hold true. The reality is that there is honor  
211 and respect amongst programmers. The manager could influence the style of junior  
212 programmers. The reason is because it is unlikely for junior programmers to start writing  
213 an entire code repository for their first industry project. In practice, the senior engineers  
214 usually lead the junior employees into the doors and get them started by some shadowing.  
215 This formality ensure a safe transition to deliver sophisticated code production but also  
216 is able to train new employees for higher-level job tasks. However, this nature implies  
217 that it is likely for the junior employees to take on some of the flavor or style from senior  
218 engineers. Since this is the first real project, junior employees may take it to heart and  
219 set it as a "benchmark" for all future project. This means that the coding styles can be  
220 passed from a senior engineer to a junior employee, which means the coding style is not  
221 unique to each individual.

222 A third assumption is the size of the code in production follows that of the dataset. This  
223 assumption cannot be held true in practice. In a clean, processed, and manipulated data  
224 set, each instance could have similar length. In practice, it is unlikely that each code  
225 project has the same length. In similar job tasks, this can sometimes be true. However,  
226 it is a very difficult assumption to defend. In reality, the length of the code tends to be  
227 dependent on the nature of the job tasks. Machine learning code can be short when  
228 the model is called from a previously trained project. The same code can also be long  
229 because sometimes it may be beneficial to design a customized model that can be tailored  
230 to the particular dataset in hand which can produce extremely lengthy code.

## 231 **References**

232 Abbasi, A. and Chen, H. (2008). Writeprints: A stylometric approach to identity-  
233 level identification and similarity detection in cyberspace. *ACM Transactions on*  
234 *Information Systems (TOIS)*, 26(2):1–29.

235 Abuhamad, M., Abuhmed, T., Mohaisen, D., and Nyang, D. (2021). Large-scale and  
236 robust code authorship identification with deep feature learning. *ACM Transactions*  
237 *on Privacy and Security (TOPS)*, 24(4):1–35.

238 Afroz, S., Islam, A. C., Stolerman, A., Greenstadt, R., and McCoy, D. (2014). Doppel-  
239 gänger finder: Taking stylometry to the underground. In *2014 IEEE Symposium on*  
240 *Security and Privacy*, pages 212–226. IEEE.

241 Burrows, S. and Tahaghoghi, S. M. (2007). Source code authorship attribution using  
242 n-grams. In *Proceedings of the twelfth Australasian document computing symposium,*  
243 *Melbourne, Australia, RMIT University*, pages 32–39. Citeseer.

244 Burrows, S., Tahaghoghi, S. M., and Zobel, J. (2007). Efficient plagiarism detection for  
245 large code repositories. *Software: Practice and Experience*, 37(2):151–175.

246 Burrows, S., Uitdenbogerd, A. L., and Turpin, A. (2014). Comparing techniques for  
247 authorship attribution of source code. *Software: Practice and Experience*, 44(1):1–32.

- 248 Caliskan-Islam, A., Harang, R., Liu, A., Narayanan, A., Voss, C., Yamaguchi, F., and  
249 Greenstadt, R. (2015). De-anonymizing programmers via code stylometry. In *24th*  
250 *USENIX security symposium (USENIX Security 15)*, pages 255–270.
- 251 Frantzeskou, G., Stamatatos, E., Gritzalis, S., and Katsikas, S. (2006). Effective identi-  
252 fication of source code authors using byte-level information. In *Proceedings of the*  
253 *28th international conference on Software engineering*, pages 893–896.
- 254 Juneja, B. (2009). *Programming with C++*. New Age International.
- 255 Juola, P. (2008). *Authorship attribution*, volume 3. Now Publishers Inc.
- 256 Kešelj, V., Peng, F., Cercone, N., and Thomas, C. (2003). N-gram-based author profiles  
257 for authorship attribution. In *Proceedings of the conference pacific association for*  
258 *computational linguistics, PACLING*, volume 3, pages 255–264.
- 259 Koppel, M., Schler, J., and Argamon, S. (2009). Computational methods in authorship  
260 attribution. *Journal of the American Society for information Science and Technology*,  
261 60(1):9–26.
- 262 Malin, C. H., Casey, E., and Aquilina, J. M. (2008). *Malware forensics: investigating*  
263 *and analyzing malicious code*. Syngress.
- 264 Sinkov, A. (1963). Programming the ibm 1401: A self-instructional programmed  
265 manual.
- 266 Spafford, E. H. and Weeber, S. A. (1993). Software forensics: Can we track code to its  
267 authors? *Computers & Security*, 12(6):585–595.
- 268 Stamatatos, E. (2009). A survey of modern authorship attribution methods. *Journal of*  
269 *the American Society for information Science and Technology*, 60(3):538–556.
- 270 Uzuner, Ö. and Katz, B. (2005). A comparative study of language models for book  
271 and author recognition. In *International conference on natural language processing*,  
272 pages 969–980. Springer.
- 273 Wilcox, L. J. (1998). Authorship: the coin of the realm, the source of complaints. *Jama*,  
274 280(3):216–217.
- 275 Wilkes, M. V., Wheeler, D. J., and Gill, S. (1951). *The Preparation of Programs for an*  
276 *Electronic Digital Computer: With special reference to the EDSAC and the Use of a*  
277 *Library of Subroutines*. Addison-Wesley Press.