

Fig. 1 Scheme of the model's training process by applying the proposed technique from the text input sequence

### A. Data Preprocessing

Before the code data is fed as an input sequence of the model, cleaning the raw text must be performed because of the characteristic of the code. The code may be written in an unorganized manner by the coder or contain text information that is irrelevant to compilation. Especially, the existing typo information in the code is unnecessary, and it can act as a fatal factor in the model training procedure. First, unnecessary new lines and spaces are removed. Comments (sequences marked with '#') written to help understand the code are also somewhat unnecessary for our training process. Hence, they are removed along with white spaces.

### B. Data Filtering Strategies

Duplication of data may occur when using different datasets as training and test data. Therefore, a model trained with these duplicated features cannot perform an accurate evaluation. Before enhancing the model with additional data, we observed data duplication between the additional training data and the test data. Hence, filtering was performed with three strategies. Algorithm 1 describes the three-step data filtering strategies in Fig. 1 more effectively.

1) *Deduplication of Data*: First, the data corresponding to A and B, and the TEST codes are loaded. The first filtering strategy filters overlapped sequence data using a hash table (HT). This method records all values of A in the HT and then checks whether these values exist in B; this is the most common and widely used method. For the first filtering of Algorithm 1 (1~6 lines), (1) the sequence input to the DEDUPLICATION function and initialized HT creates a new HT from the A, and (2) the new HT is compared with B to create *First\_filtered\_codes* (6 lines). Most duplicate data are filtered out in this process.

2) *Deletion of Intersection*: Second, the purpose of the *Second Filtering* is to filter out data not completely filtered out by the HT owing to reasons such as trailing space. This includes spaces such as white spaces ' ' in the right or left edge of sentences, tabs '\t', In the second filtering of Algorithm 1 (7~10 lines), (1) the SIMPLIFY function concatenates all newlines existing in the character strings of the code. (2) All spaces and newlines before and after the character strings are removed. (3) Filtering is performed once more by taking the intersection of these filtered sequences of the test code and using *first\_filtered\_codes* to generate *second\_filtered\_codes*.

3) *Exhaustive Search*: Most of the duplicated data are removed by the second filtering process. However, a method that completely eliminates duplication is required. Therefore, an exhaustive search is performed, and s and *TEST\_codes* are mutually compared in the third filtering of Algorithm 1 to remove the few remaining duplicate data traces (11~20 lines).

Now, the *third filtered codes*, which are used as the training data, are finally generated.

### C. Forming positive and negative pairs

Then, the dataset that has been trimmed with the previous three preprocessing strategies should be composed of positive and negative pairs for the training process. One type of training data is positive pairs which can determine whether the two codes are similar, and the other is negative pairs which can determine whether the two codes are not similar. For negative samples, ranks are assigned after sorting the data in descending order based on the confidence score. BM25 algorithm of the Okapi system, which is based on probabilistic retrieval research, is a ranking function utilized by search engines to rank matching sentences according to their relevance to a given query. For length normalization of BM25, BM25L, a newer variant to boost scores of very long documents and more effective than BM25, was proposed [25]. The Okapi BM25 scores a document  $D$  with respect to query sentence  $Q$  containing keywords  $k_1, \dots, k_n$  as follows:

$$score(D, Q) = \sum_{i=1}^n IDF(k_i) \cdot \frac{f(k_i, D) \cdot (m_1 + 1)}{f(k_i, D) + m_1 \cdot (b \cdot \frac{|D|}{adl} + 1 - b)}$$

where  $f(k_1, D)$  is the term frequency of  $k_1$  in  $D$ ,  $|D|$  is the length of the words, and  $adl$  is average document length.  $m_1$  and  $b$  are hyperparameters.

$$IDF(k_i) = \ln\left(\frac{N - n(k_i) + 0.5}{n(k_i) + 0.5}\right) + 1$$

The inverse document frequency (IDF) means how common or rare a word is in the total document set.  $IDF(k_i)$  is consisted of IDF weight of the query term  $k_i$  and  $n(k_i)$  is the number of documents containing  $k_i$ , which is calculated as above. Using these algorithms, the volume of the original dataset was dramatically expanded, and we experimented with demonstrating the improved efficiency of the code similarity evaluation model.

Algorithm 1 Efficient Data Filtering Strategies
A = DATA1, B = DATA2
Load TEST_codes
<b>procedure</b> DEDUPLICATION(sequence, HT)
→ Quick check if $i$ is in HT
1: Read All A, B values
2: <b>for</b> $i = 0$ to $\text{len}(A[i])$ <b>do</b>
3: HT.add(A[i])
→ Build HT
4: <b>for</b> $i = 0$ to $\text{len}(A[i])$ <b>do</b>
5: <b>if</b> B[i] not in HT <b>then</b>
6: First_filtered_codes.append(B[i])
→ First Filtering
<b>procedure</b> SIMPLIFY(code)
7: <b>return</b> <code>','.join(code.split('\n')).rstrip(' ').strip()</code>

```

procedure DELETE INTERSECTION(First_filtered_codes)
8: for  $i = 0$  to  $\text{len}(\text{First\_filtered\_codes}[i])$  do
9:   if SIMPLIFY(First_filtered_codes[i]) not in
INTERSECTION(TEST_codes) then
10:
Second_filtered_codes.append(First_filtered_codes[i])
→ Second Filtering
procedure EXHAUSTIVE SEARCH(Second_filtered_codes)

11: for  $i = 0$  to  $\text{len}(\text{Second\_filtered\_codes}[i])$  do
12:   USE_TOKEN = True
13:   if Second_filtered_codes[i] in TEST_codes then
14:     continue
15:   else
16:     for  $s = 0$  to TEST_codes do  $\text{len}(s) > 0$  and
 $\text{len}(\text{Second\_filtered\_codes}[i]) >$  and  $((\text{Second\_filtered\_codes}[i]$ 
in  $s)$ 
           or  $(s$  in  $\text{Second\_filtered\_codes}[i])$ )
17:     USE_TOKEN = False
18:     if USE_TOKEN == True then
19:
Third_filtered_codes.append(Second_filtered_codes[i])
→ Third Filtering
20: Return Third filtered codes

```

#### D. Cross-Validation Ensemble

The training process of the language models is divided into training, validating, and testing datasets, respectively. For evaluation of our model, any data imbalance-related issues can not be allowed; if there are, an accurate measure is impossible. In Fig. 2, the 5-fold cross-validation can train each model in the training step, alternating each set 5 times. We finally divided the overall training dataset by 8:2, finetuned the model with each train & valid features, and evaluated it on the test dataset.

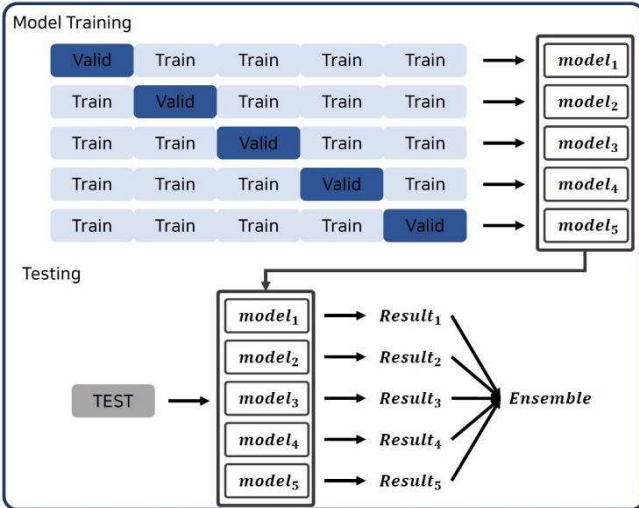


Fig. 2 A method of 5-fold cross-validation. The models were ensemble by voting from each  $model_{i(n=1,2,3,4,5)}$  and resulting in these models' training and testing

Consequently, a more accurate performance evaluation can be performed. We utilized the method of evaluating the final performance by recording the checkpoint that indicates the best performance for each fold and creating an ensemble of the determinations made by the 5 models generated in this method in the testing process.

### III. RESULTS AND DISCUSSION

This section describes the quantitative experimental method and experimental results that prove the effectiveness of our proposed method.

#### A. Data

The experiment was conducted by randomly extracting only 50% of the CodeNet data, owing to resource problems, presented at a code similarity competition, which was a NLP AI contest held in Korea in May 2022. The problems were numbered, and there were two solution codes (1 and 2) for solving each problem.

TABLE I  
VOLUME AND COUNT OF DATA FOR COMPETITION AND CODENET. ONLY 50% OF CODENET WAS USED OWING TO A RESOURCE ISSUE. THE POSITIVE AND NEGATIVE PAIRS CREATED IN THIS PROCESS ARE COMBINED AND USED AS THE TOTAL TRAINING DATA.

DATA	Competition	CodeNet
Count	45,000	120,000
Volume	185.5MB	493.25MB
Positive	Total Train:	Train: 10,316,052
& Count	5,445,594	Valid: 117,168
Negative	Valid: 63,242	
Pairs	Total Train: 3.52GB	Train: 7.34GB
	Valid: 55.1MB	Valid: 117.4MB

Table 1 shows the count and volume of the competition and CodeNet data used in the experiment. Positive and negative pairs were created to increase the volume of the data. In addition, the validation loss was calculated every 1000 steps by classifying the data as training and validation data. When two codes are input, the accuracy is used to measure the performance of the binary classification task, which determines whether the two codes solve the same problem.

#### E. Experiment Setup

The experiment environment is as follows. The graphcodebert and codebert-mlm models were used for the proposed method. Both models were trained on a large dataset used for researching the probing capability of code-based PLMs. These PLMs can learn the general-purpose representation as well as code search and documentation generation tasks. Moreover, these PLMs are language models, and studies have been recently conducted on these models for downstream tasks related to code as in Table 2.

TABLE II  
EXPERIMENTAL RESULTS OF THE LANGUAGE MODEL TRAINED ON THE KOREAN NUMERICAL REASONING DATA

Models	Competition Only		Competition + CodeNet	
	BM25	BM25L	BM25	BM25L
	Accuracy			
graphcodebert	95.7206	96.3735	96.4251	97.3456
codebert-mlm	94.9471	95.8456	95.6421	96.8657
Cross validated ensemble	96.0823	97.0988	97.6552	98.8785

Meanwhile, Table 3 lists the hyperparameters for finetuning the two models using the proposed method. These models were trained using an A100 GPU and by setting the learning rate to  $2e-5$ , epsilon to  $1e-5$ , the optimizer to AdamW, max epochs to 2, batch size to 32, and max sequence length to 512.

TABLE III  
HYPERPARAMETER SETTINGS FOR THE FINETUNING PROCESS WITH  
GRAPHCODEBERT AND CODEBERT-MLM.

Hyperparameters	Models
Learning rate	2e-5
Epsilon	1e-5
Optimizer	AdamW
Max epochs	2
Batch size	32
Max sequence length	512

#### F. Results

Table 3 shows the quantitative results of the code similarity evaluation model. The first column shows the single models used in the experiment, graphcodebert and codebert-mlm, and cross-validation ensemble of the two models. In addition, the competition-only and competition + CodeNet data were used as the training data. The similarity was measured by applying the BM25 and BM25L similarity algorithms of the Okapi system to each dataset to generate the positive and negative pairs.

The results show a consistent performance improvement when the BM25L algorithm was applied, unlike with the BM25 algorithm. In addition, the graphcodebert model (the first row of Models) outperformed the codebert-mlm model (the second row of Models). Moreover, the cross-validation ensemble of the two models achieved a better performance than the two single models. In particular, the cross-validation ensemble model finetuned with the competition + CodeNet dataset achieved the highest score with a performance of 98.8785.

#### IV. CONCLUSION

In this study, we proposed an effective code similarity representation model. Our approach of data filtering strategies: deduplication, delete intersection, and exhaustive search, can eliminate data redundancy and contributes to efficient model training. Subsequently, the positive and negative pairs were formed by applying the BM25 and BM25L algorithms to the filtered dataset. Then, the model was finetuned with the 5-fold cross-validation method. Finally, the best performance was achieved by performing a final inference on the test dataset by the five models, each of which was trained with the cross-validated ensemble method. The performance improvement was demonstrated quantitatively by comparing our approach with conventional methods. In addition, the ensemble was proven effective on both single and augmented data. Our approach to improving the productivity of this software is crucial as the impact of labor-supply shortages in IT can negatively affect both productivity and innovation, and we believe our methods will be leveraged in various NLP tasks.

#### ACKNOWLEDGMENTS

Hanshin University Research Grant supported this work.

#### REFERENCES

- [1] D. Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).
- [2] A. Radford, et al. "Language models are unsupervised multitask learners." *OpenAI blog* 1.8 (2019): 9.
- [3] Y. Zhilin, et al. "Xlnet: Generalized autoregressive pre-training for language understanding." *Advances in neural information processing systems* 32 (2019).
- [4] D. Zihang, et al. "Transformer-xl: Attentive language models beyond a fixed-length context." *arXiv preprint arXiv:1901.02860* (2019).
- [5] A. Vaswani, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).
- [6] G. Kim, et al. "AI Student: A Machine Reading Comprehension System for the Korean College Scholastic Ability Test." *Mathematics* 10.9 (2022): 1486.
- [7] S. Lee, G. Kim, and H. Lim, "Verification of educational goal of reading area in Korean SAT through natural language processing techniques," *Journal of the Korea Convergence Society*, vol. 13, no. 1, pp. 81–88, Jan. 2022.
- [8] G. Kim, et al. "Automatic extraction of named entities of cyber threats using a deep Bi-LSTM-CRF network." *International journal of machine learning and cybernetics* 11.10 (2020): 2341-2355.
- [9] K. Kim, et al. "GREG: A global level relation extraction with knowledge graph embedding." *Applied Sciences* 10.3 (2020): 1181.
- [10] Z. Feng, et al. 2020. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1536–1547, Online. Association for Computational Linguistics.
- [11] G. Daya, et al. "Graphcodebert: Pre-training code representations with data flow." *International Conference on Learning Representations: ICLR 2021*.
- [12] C. Kevin, et al. "Electra: Pre-training text encoders as discriminators rather than generators." *arXiv preprint arXiv:2003.10555* (2020).
- [13] T. Breaux and J. Moritz. 2021. The 2021 software developer shortage is coming. *Commun. ACM* 64, 7 (July 2021), 39–41. <https://doi.org/10.1145/3440753>.
- [14] P. Tambe, Xuan Ye, Peter Cappelli (2020) Paying to Program? Engineering Brand and High-Tech Wages. *Management Science* 66(7):3010-3028. <https://doi.org/10.1287/mnsc.2019.3343>.
- [15] H. Hamel, et al. "Codesearchnet challenge: Evaluating the state of semantic code search." *arXiv preprint arXiv:1909.09436* (2019).
- [16] L. Yujia, et al. "Competition-level code generation with alphacode." *arXiv preprint arXiv:2203.07814* (2022).
- [17] Z. Feng, et al. "Flowchart-based cross-language source code similarity detection." *Scientific Programming* (2020).
- [18] S. Ducasse, M. Rieger and S. Demeyer, "A language independent approach for detecting duplicated code," *Proceedings IEEE International Conference on Software Maintenance - 1999 (ICSM'99)*. 'Software Maintenance for Business Change' (Cat. No.99CB36360), 1999, pp. 109-118, doi: 10.1109/ICSM.1999.792593.
- [19] B. S. Baker, "On finding duplication and near-duplication in large software systems," *Proceedings of 2nd Working Conference on Reverse Engineering*, 1995, pp. 86-95, doi: 10.1109/WCRE.1995.514697.
- [20] I. D. Baxter, A. Yahin, L. Moura, M. Sant'Anna and L. Bier, "Clone detection using abstract syntax trees," *Proceedings. International Conference on Software Maintenance (Cat. No. 98CB36272)*, 1998, pp. 368-377, doi: 10.1109/ICSM.1998.738528.
- [21] M. Leblanc and Merlo, "Experiment on the automatic detection of function clones in a software system using metrics," 1996 *Proceedings of International Conference on Software Maintenance*, 1996, pp. 244-253, doi: 10.1109/ICSM.1996.565012.
- [22] P. Anupriya. *Code Clone Detection Using Code2Vec*. University of California, Irvine, 2020.
- [23] A. Uri, et al. "code2vec: Learning distributed representations of code." *Proceedings of the ACM on Programming Languages* 3.POPL (2019): 1-29.
- [24] R. Stephen, et al. "Okapi at TREC-3." *Nist Special Publication Sp 109* (1995): 109.
- [25] Y. Lv and C. Zhai. 2011. When documents are very long, BM25 fails! In *Proceedings of the 34th international ACM SIGIR conference on research and development in Information Retrieval (SIGIR '11)*. Association for Computing Machinery, New York, NY, USA, 1103–1104. <https://doi.org/10.1145/2009916.2010070>.
- [26] K. Yang, et al. "Cross-Validated Ensemble Methods in Natural Language Inference." *Annual Conference on Human and Language Technology*. Human and Language Technology, 2019.