

Appendix

1 Evaluation Metrics

We provide the details of the evaluation metrics we used in the experiment.

1.1 BLEU

BLEU measures the average n-gram precision between the reference sentences and generated sentences, with a brevity penalty for short sentences. The formula to compute BLEU-1/2/3/4 is:

$$\text{BLEU-N} = BP \cdot \exp \sum_{n=1}^N \omega_n \log p_n, \quad (1)$$

where p_n (n-gram precision) is the fraction of n-grams in the generated sentences which are present in the reference sentences, and ω_n is the uniform weight $1/N$. Since the generated summary is very short, high-order n-grams may not overlap. We use the +1 smoothing function (Lin and Och, 2004). BP is a brevity penalty given as:

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases} \quad (2)$$

Here, c is the length of the generated summary, and r is the length of the reference sentence.

1.2 Rouge-L

Based on the longest common subsequence (LCS), Rouge-L is widely used in text summarization. Instead of using only recall, it uses the F-score which is the harmonic mean of precision and recall values. Suppose A and B are generated and reference summaries of lengths c and r respectively, we have:

$$\begin{cases} P_{\text{Rouge-L}} = \frac{\text{LCS}(A,B)}{\text{LCS}(A,B)} \\ R_{\text{Rouge-L}} = \frac{\text{LCS}(A,B)}{r} \end{cases} \quad (3)$$

$F_{\text{Rouge-L}}$, which indicates the value of Rouge-L, is calculated as the weighted harmonic

mean of $P_{\text{Rouge-L}}$ and $R_{\text{Rouge-L}}$:

$$F_{\text{Rouge-L}} = \frac{(1 + \beta^2) P_{\text{Rouge-L}} \cdot R_{\text{Rouge-L}}}{R_{\text{Rouge-L}} + \beta^2 P_{\text{Rouge-L}}} \quad (4)$$

β is set to 1.2 as in (Zhang et al., 2020; Wan et al., 2018).

1.3 Meteor

Meteor is a recall-oriented metric that measures how well the model captures the content from the references in the generated sentences and has a better correlation with human judgment. Suppose m is the number of mapped unigrams between the reference and generated sentence with lengths c and r respectively. Then, precision, recall, and F are given as:

$$P = \frac{m}{c}, \quad R = \frac{m}{r}, \quad F = \frac{PR}{\alpha P + (1 - \alpha)R} \quad (5)$$

The sequence of mapping unigrams between the two sentences is divided into the fewest possible number of ‘‘chunks’’. This way, the matching unigrams in each ‘‘chunk’’ are adjacent (in two sentences) and the word order is the same. The penalty is then computed as:

$$\text{Pen} = \gamma \cdot \text{frag}^\beta \quad (6)$$

where frag is a fragmentation fraction: $\text{frag} = ch/m$, where ch is the number of matching chunks and m is the total number of matches. The default values of α, β, γ are 0.9, 3.0 and 0.5 respectively.

1.4 Cider

Cider is a consensus-based evaluation metric used in image captioning tasks. The notions of importance and accuracy are inherently captured by computing the TF-IDF weight for each n-gram and using cosine similarity for sentence similarity. To compute Cider, we first calculate the TF-IDF weighting $g_k(s_i)$ for each n-gram ω_k in reference

sentence s_i . Here ω is the vocabulary of all n-grams. Then we use the cosine similarity between the generated sentence and the reference sentences to compute Cider_n score for n-grams of length n . The formula is given as:

$$\text{Cider}_n(c_i, s_i) = \frac{\langle \mathbf{g}^n(c_i), \mathbf{g}^n(s_i) \rangle}{\|\mathbf{g}^n(c_i)\| \|\mathbf{g}^n(s_i)\|} \quad (7)$$

where $\mathbf{g}^n(s_i)$ is a vector formed by $g_k(s_i)$ corresponding to all the n-grams (n varying from 1 to 4). c_i is the i^{th} generated sentence. Finally, the scores of various n-grams can be combined to calculate Cider as follows:

$$\text{Cider}(c_i, s_i) = \sum_{n=1}^N w_n \text{Cider}_n(c_i, s_i) \quad (8)$$

2 Applying our Framework to Seq2Seq Neural Network Models

We adapt our framework to four neural-based models, NMTGen (M1), CommitBERT (M2), CodeT5-small (M3) and CodeT5-base (M4). The experimental results in different metrics are shown in Figure 1, 2, 3, and 4. We can see that our framework can improve the performance of all four neural models on all programming languages (except CommitBERT on JavaScript in terms of Meteor). These results demonstrate that our approach can be applied to these neural models and boost their performance.

3 Human Evaluation

We conduct a human evaluation to evaluate the effectiveness of the commit message generated by our approach RACE and the other four approaches (NNGen, NMTGen, CommitBERT, and CoRec) in terms of three aspects, Informative, Conciseness and Expressiveness. We define detailed criteria on three aspects for manual labeling in Table 1.

The experimental results in the main text show that RACE outperforms the other approaches in three aspects. We also confirmed the dominance of our approach using Wilcoxon signed-rank tests for human evaluation. The results shown in Table 2 reflect that the improvement of RACE over other approaches is statistically significant with all p-values smaller than 0.05 at 95% confidence level.

References

Chin-Yew Lin and Franz Josef Och. 2004. ORANGE: a method for evaluating automatic evaluation metrics for machine translation. In *COLING*.

Yao Wan, Zhou Zhao, Min Yang, Guandong Xu, Haochao Ying, Jian Wu, and Philip S. Yu. 2018. Improving automatic source code summarization via deep reinforcement learning. In *ASE*.

Jian Zhang, Xu Wang, Hongyu Zhang, Hailong Sun, and Xudong Liu. 2020. Retrieval-based neural source code summarization. In *ICSE*.

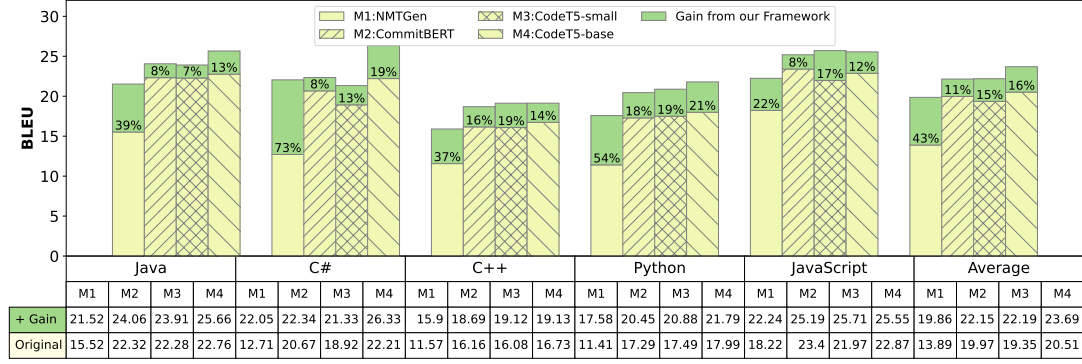


Figure 1: Performance gains on four models. The results are evaluated under BLEU.

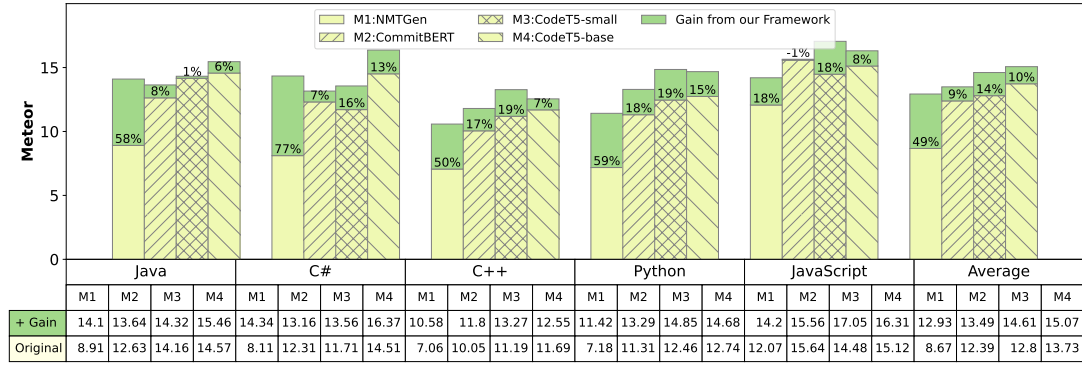


Figure 2: Performance gains on four models. The results are evaluated under Meteor

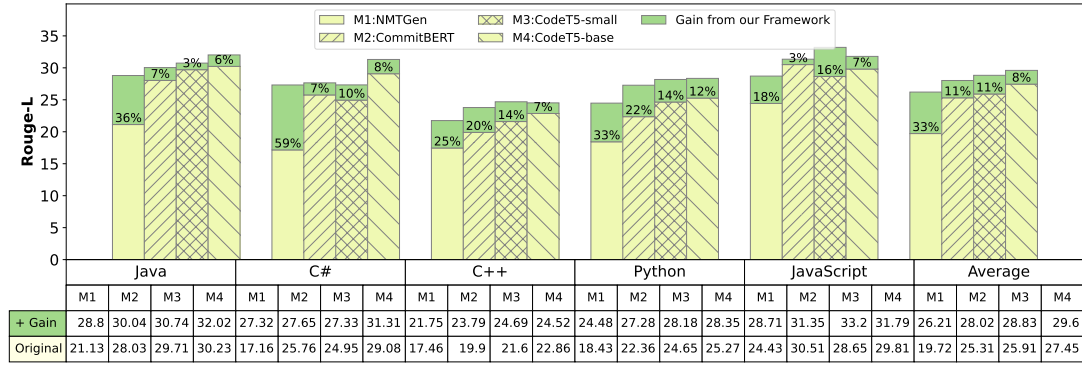


Figure 3: Performance gains on four models. The results are evaluated under Rouge-L.

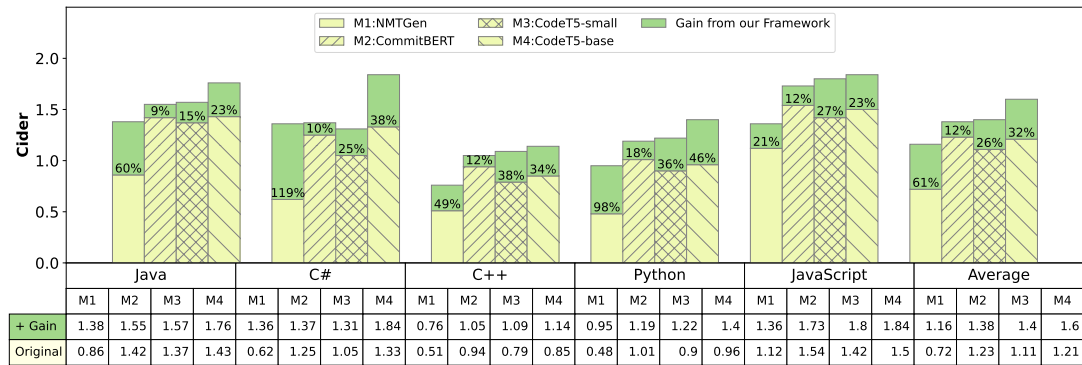


Figure 4: Performance gains on four models. The results are evaluated under Cider.

Score	Meaning
<u>Informativeness</u>	
Is the important information about the code diff reflected in the commit message?	
0	Missing all of the information about the code diff.
1	Is missing some very important information that can hinder the understanding of the code diff.
2	Is missing some information but some of the missing information is not necessary to understand the code diff.
3	Is missing some information but all of the missing information is not necessary to understand the code diff.
4	Is not missing any information.
<u>Conciseness</u>	
Is there extraneous information included in the commit message?	
0	All of the information is unnecessary.
1	Has a lot of unnecessary information.
2	Has some unnecessary information.
3	Has a little unnecessary information.
4	Has no unnecessary information.
<u>Expressiveness</u>	
How readable and understandable is the commit message?	
0	Cannot read and understand.
1	Is hard to read and understand.
2	Is somewhat readable and understandable.
3	Is mostly readable and understandable.
4	Is easy to read and understand.

Table 1: The meaning of scores in human evaluation on three metrics.

Model	Informativeness	Conciseness	Expressiveness
CommitBERT	$1.21e^{-24}$	$3.12e^{-22}$	$7.90e^{-05}$
NNGen	$2.04e^{-27}$	$2.28e^{-25}$	$4.28e^{-08}$
NMTGen	$8.64e^{-30}$	$5.87e^{-28}$	$2.55e^{-13}$
CoRec	$6.90e^{-26}$	$2.57e^{-26}$	$1.27e^{-06}$

Table 2: Statistics significance p-value of RACE over other approaches in human evaluation.