

CAMCMG: A Change-aware Model for Commit Message Generation

Anonymous ACL submission

Abstract

Commit messages play a crucial role in understanding code change intent and tracking software evolution, and automatically generating high-quality commit messages is important for improving software maintenance efficiency. Although existing methods have made preliminary attempts to capture the semantics of *diff*, semantic associations across hunks and the semantic alignment between the *diff* and its commit message have not been further explored. In this paper, we propose a change-aware model for commit message generation called CAMCMG, which focuses on the hunk associations in code change and alignment between code change and commit message. Specifically, it adopts a local attention mechanism based on a sliding window to capture local contextual information around hunks, a CMG-oriented change attention mechanism to model semantic associations across hunks, and a change alignment loss to optimize the model to enhance the generation of change-aligned tokens during the training phase. Experimental results show that CAMCMG achieves an average improvement of 6.4% on automatic evaluation metrics and 5.2% on human evaluation metrics. Ablation studies further demonstrate the effectiveness of each component, where local attention, change attention, and change alignment loss contribute average improvements of 5.1%, 9.6%, and 7.1%, respectively, on the automatic evaluation metrics.

1 Introduction

Commit message summarizes the purpose and intent of code change¹, which plays an important role in program comprehension, such as code review (Caulo et al., 2020; Li et al., 2025), defect localization (Wang et al., 2025; Li et al., 2021; Schlüter et al., 2021) and software evolution analysis (Gîrba and Ducasse, 2006; Wei et al., 2025).

¹In this paper, we use *diff* and *code change* interchangeably.

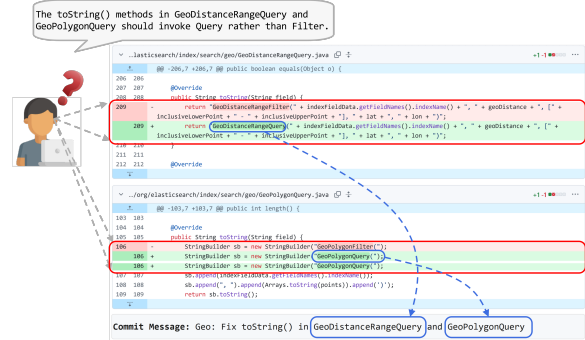


Figure 1: Two hunks modify the code by changing the suffix of GeoPolygonFilter and GeoDistanceRangeFilter to Query. The commit message summarizes the modified methods and classes reflected in the diff.

High-quality commit messages can help developers quickly understand the context of changes and improve the efficiency of team collaboration. However, there are many commit messages that are vague or even blank in real-world projects (Maalej and Happel, 2010), since writing commit messages manually is time-consuming and tedious (Zhang et al., 2024c). For example, approximately 14% of the commit messages in over 23,000 SourceForge projects were completely empty (Dyer et al., 2013); and this problem also existed in the Apache Software Foundation projects (Imani et al., 2024). Therefore, various automatic commit message generation (CMG) methods have been proposed to generate accurate, concise, and semantically logical descriptions for code changes. Specifically, given a code change $diff = \{h_1, h_2, \dots, h_n\}$, where h_i denotes the i -th hunk, CMG aims to transform *diff* to a commit message.

Existing CMG methods can be broadly classified into template-based (Buse and Weimer, 2010; Shen et al., 2016), retrieval-based (Huang et al., 2020; Liu et al., 2018; Hoang et al., 2020), learning-based (Dong et al., 2022; Liu et al., 2023; Nie

et al., 2021; Wang et al., 2024a) and hybrid methods (Liu et al., 2020; Shi et al., 2022; Zhang et al., 2024b; He et al., 2023). Recently, with the rising of pre-trained language models (PLMs), PLM-based methods become dominant in learning-based methods. These methods are based on the Transformer architecture, which adopts the standard full attention mechanism (Vaswani et al., 2017). However, standard full attention mechanism computes the attention scores among arbitrary tokens, resulting in the attention being diffusely distributed (Correia et al., 2019). Many studies attempt to leverage a sparse attention mechanism to mitigate this issue (Beltagy et al., 2020; Guo et al., 2023; Wang et al., 2024b). Despite their efforts, there is still a gap between the CMG scenario and the general scenario. The following two key observations provide an opportunity to design a CMG-oriented sparse attention mechanism.

(O1): Semantic associations among hunks. Even if the different hunks in a *diff* are positionally separated, they are not completely semantically independent. As shown in Figure 1, the *diff* modifies the `toString()` method. Although the two hunks are separate in location, they share the same intent, i.e., calling class `Query` instead of `Filter`. Therefore, modeling cross-hunk semantic associations is necessary to capture the intent of code changes.

(O2): Semantic alignment between *diff* and commit message. Tokens in the commit message often correspond to certain changes in the *diff* (e.g., function calls, variable usage). As shown in Figure 1, the commit message summarizes the changes in *diff*, which involve the `toString()` method in both the "GeoDistanceRangeQuery" and "GeoPolygonQuery" classes. It emphasizes the importance of aligning the semantic content between the *diff* and its commit message.

In summary, generating high-quality commit messages requires modeling semantic associations across multiple hunks and aligning the generated message with the *diff*. However, to the best of our knowledge, *off-the-peg* methods do not sufficiently satisfy these *desiderata*. Therefore, this leads to the major challenges: how to associate the semantics across hunks and how to align the *diff* and its corresponding commit message.

Motivated by this, we propose CAMCMG, which is a Change-aware Model for Commit Message Generation. First, context plays an important role in commit message generation because the effectiveness of the generated message de-

pends on the surrounding context, except for code changes (Sillito et al., 2008). Therefore, we employ a sliding window to capture local contextual information around each hunk. It allows each token to pay attention to the neighboring tokens. Second, we further introduce a CMG-oriented change attention mechanism to capture semantic associations across hunks, which enables related tokens among different hunks to establish attention connections. Finally, we design a change alignment loss to optimize the model to enhance the generation of tokens in the commit message that correspond to the changes of *diff* (denoted as change-aligned tokens) during the training phase.

To evaluate the effectiveness of CAMCMG, we conduct extensive experiments on MCMD, one of the most widely used datasets for CMG, and compare it with ten state-of-the-art (SOTA) methods. Experimental results show that CAMCMG achieves the best overall performance, with an average improvement of 6.4% on automatic evaluation metrics and 5.2% on human evaluation metrics compared with baseline models except for GPT-4. Ablation studies further confirm the effectiveness of each component. Specifically, local attention, change attention, and change alignment loss contribute average improvements of 5.1%, 9.6%, and 7.1%, respectively, on the automatic evaluation metrics. The main contributions of our paper can be summarized as follows:

- We propose CAMCMG, which includes local attention and change attention. Specifically, *local attention* focuses on capturing the local contextual information around each hunk, while *change attention* is designed specifically for the CMG task to model the semantic associations across hunks.
- We design a change alignment loss to optimize the model to enhance the generation of change-aligned tokens during the training phase, which directs the model to pay more attention to the semantic alignment tokens, thus improving the relevance and accuracy of commit messages.
- Extensive experiments demonstrate that CAMCMG outperforms the state-of-the-art methods on multiple evaluation metrics, which show its effectiveness and superiority.

2 Related Work

2.1 Commit Message Generation

Researchers have proposed various methods to generate commit messages. Specifically, it can be categorized as templated-based, retrieval-based, learning-based, hybrid and LLM-based methods. Early work generated commit messages based on predefined templates (Buse and Weimer, 2010; Shen et al., 2016). For example, Buse et al. (Buse and Weimer, 2010) leveraged control flow patterns to generate commit messages. Retrieval-based methods retrieved similar commits from the training set and reused their commit messages (Liu et al., 2018; Huang et al., 2020; Hoang et al., 2020). For example, Liu et al. (Liu et al., 2018) used cosine similarity to retrieve similar code changes. Learning-based methods adopted neural network models to translate *diffs* into commit messages. For example, Dong et al. (Dong et al., 2022) proposed FIRA, which builds fine-grained graphs by combining sub-tokens nodes with edit operation nodes to capture small differences in code changes. Liu et al. (Liu et al., 2023) proposed CCRep, which uses a pre-trained model with a query-back mechanism to model interactions between code changes and their contexts. Hybrid methods leverage both retrieval-based and learning-based approaches to generate commit messages. For example, ATOM (Liu et al., 2020) extracts semantic information from code changes using ASTs and combines retrieval and generation through hybrid ranking. RACE (Shi et al., 2022) retrieves similar code-message pairs to guide generation. COME (He et al., 2023) learns contextualized code change representations via modification embeddings and a self-supervised objective, and integrates retrieval and generation through a decision algorithm. Recent work has explored LLMs for commit message generation. Wu et al. (Wu et al., 2025) empirically validated the effectiveness and generalizability of LLMs through in-context learning.

Although the above methods have achieved promising performance, the semantic alignment between *diffs* and commit messages is rarely modeled explicitly, which potentially affects the accuracy and expressiveness of the resulting messages.

2.2 Sparse Attention Mechanisms

The self-attention mechanism allows each token to attend to all others (Zaheer et al., 2020). In long code sequences, this indiscriminate attention can

dilute focus on semantically related tokens across distant hunks, reducing the model’s ability to capture code semantics. To mitigate this, prior work has explored sparse attention mechanisms. For example, Wang et al. (Wang et al., 2024b) proposed SparseCoder, which modeled short-term dependencies using a sliding-window mechanism and introduced global attention and identifier attention to capture long-range dependencies among source code identifiers. Guo et al. (Guo et al., 2023) proposed LongCoder, which introduced bridge attention and global attention, effectively captured long-range dependency while maintaining fine-grained local semantic representations. Zhang et al. (Zhang et al., 2024a) proposed VulAdvisor, which incorporated local attention that used TF-IDF weights to emphasize locally important code elements and introduced repair action loss that prioritized semantically meaningful tokens during training, improving the actionability and precision of the suggestions.

Although existing sparse attention methods effectively address the challenges of long sequences and capture long-range dependencies in code, they generally lack explicit modeling of the semantic relationships among different hunks within a *diff*. Consequently, cross-hunk associations that are crucial for representing the overall semantics of code changes cannot be adequately captured.

3 Methodology

The overall framework of CAMCMG is illustrated in Figure 2. Given a *diff* as input, the first step is to construct specialized mask matrices for local attention and change attention. These mask matrices serve to indicate the key regions the model should focus on, enabling the model to focus on both the local context of each hunk and the semantic associations across multiple hunks. Once the mask matrices are prepared, they are fused by taking the element-wise maximum of them. Subsequently, it is integrated into the Transformer’s self-attention mechanism. This allows the model to capture the local contextual information around hunks and associations among hunks. Formally, the standard self-attention is computed as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \quad (1)$$

$$Q = XW^Q \quad K = XW^K \quad V = XW^V, \quad (2)$$

where X is the embedded representation of the input sequence, and W^Q , W^K , and W^V are learn-

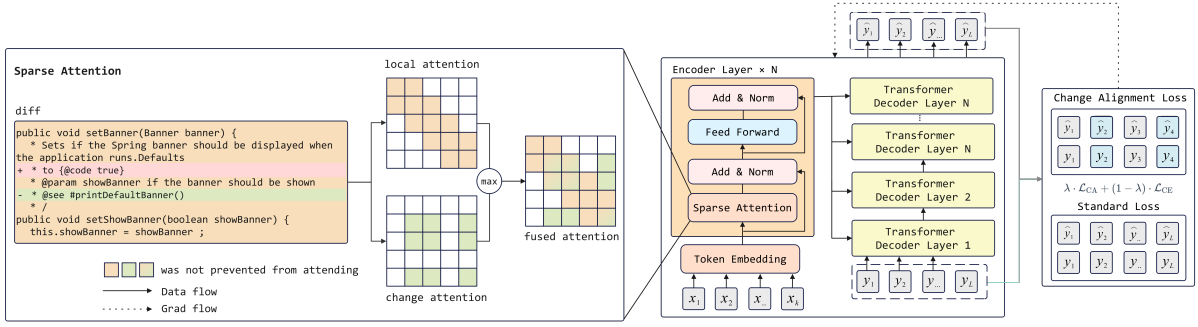


Figure 2: The overall workflow of CAMCMG.

able weight matrices. With the mask matrix incorporated, the attention computation is modified as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}} + M\right)V \quad (3)$$

where M denotes the mask matrix, which restricts the range of attention computation, thereby controlling the positions each token can attend to when constructing contextual representations. If the i -th token is allowed to attend to the j -th token, $M_{ij} = 0$; otherwise $M_{ij} = -\infty$. When M_{ij} is set to $-\infty$, the attention score from the i -th to j -th token will be zero after using the softmax function.

During the training phase, the model is further optimized with a change alignment loss, which encourages the generated message to semantically align with the key areas captured in the *diff*.

3.1 Local Attention

To enable the model to effectively capture the local contextual information around each hunk, we introduce a local attention mechanism. Unlike the self-attention mechanism in Transformers, local attention does not allow each token to attend to all others; instead, it restricts the attention scope to a limited neighborhood. This design enables the model to learn fine-grained contextual dependencies while reducing the influence of irrelevant regions. Given a local window size ω , we define a range of local attention for each token in the *diff*, which only includes $\lfloor \omega/2 \rfloor$ tokens before and after the current token. The local attention mask M^{local} is thus constructed as:

$$M_{ij}^{local} = \begin{cases} 0 & \text{if } |i - j| \leq \lfloor \omega/2 \rfloor \\ -\infty & \text{otherwise,} \end{cases} \quad (4)$$

where i and j denote the positions of tokens in the *diff* sequence.

3.2 Change Attention

The *diff* typically consists of multiple hunks, each describing a localized semantic modification. However, real-world commits often involve multiple interrelated changes that may occur in different functions, classes, or modules, leading to semantic associations that span across hunks. Although local attention can effectively capture contextual information around a hunk, it fails to model the semantic linkages across hunks, which are essential for understanding the overall intent of a code change. For example, as shown in Figure 1, the `toString()` methods in two different files were modified, where the class names in the output strings were changed from "GeoDistanceRangeFilter" and "GeoPolygonFilter" to their corresponding "GeoDistanceRangeQuery" and "GeoPolygonQuery".

Therefore, we introduce change attention, which explicitly captures semantic relations across different hunks. We identify all hunks in the *diff* and represent their token position sets as \mathcal{C} . Tokens from different hunks can attend to each other within \mathcal{C} , allowing semantic information to flow across hunks. This enables the model to aggregate related modification cues and form a unified understanding of the overall code change. The change attention mask M^{change} is thus constructed as:

$$M_{ij}^{change} = \begin{cases} 0, & \text{if } i, j \in \mathcal{C} \\ -\infty, & \text{otherwise.} \end{cases} \quad (5)$$

The above two attention mechanisms enable CAMCMG to capture and integrate semantic associations around and across hunks, constructing a more comprehensive representation of the entire code change.

The mask matrix M in Eq.(1) is defined as:

$$M = \max \left(M_{ij}^{\text{local}}, M_{ij}^{\text{change}} \right), \quad (6)$$

where $\max(\cdot, \cdot)$ is the element-wise maximum function.

3.3 Change Alignment Loss

In practice, certain tokens in the diff (such as variable names, function names, or class names) often carry crucial semantic information that is explicitly reflected in the commit message. For example, when developers rename a function or variable, or add a new method, the commit message typically describes these specific changes. This indicates that there exists a semantic alignment between tokens in the diff and the commit message, which is important for accurately generating commit messages. Therefore, during model optimization, we introduce a change alignment loss, which builds upon the standard cross-entropy loss to further enhance the model’s focus on the actual modifications in the diff. The model is able to more accurately capture the core modification information and achieve semantic alignment between the diff and its commit message by applying additional optimization constraints to the change-aligned tokens.

Algorithm 1: Change Alignment Mask Generation

Input: Change code *diff*, target commit message cm_t , generated message tokens t_{cm_g}

Output: Change alignment mask ca_mask .

- 1 $W_{diff} \leftarrow$ tokenize *diff* by spaces and lemmatize them
 - 2 $W_{cm_t} \leftarrow$ tokenize cm_t by spaces and lemmatize them
 - 3 $W_{co} \leftarrow W_{diff} \cap W_{cm_t}$
 - 4 $cm_g = \text{decode}(t_{cm_g}) \triangleright$ decode tokens to word sequence
 - 5 $ca_mask \leftarrow [0] \times |t_{cm_g}|$
 - 6 **for** w **in** cm_g **do**
 - 7 $w \leftarrow$ lemmatize w
 - 8 **if** $w \in W_{co}$ **then**
 - 9 set the elements in ca_mask whose index is equal to the index of w in t_{cm_g} to 1
 - 10 **Return** ca_mask
-

Algorithm 1 outlines the process of identifying the change-aligned tokens and creating their masks. The algorithm first tokenizes and lemmatizes both the diff and the target commit message to obtain their respective word sets (Lines 1–2). In Line 3, it identifies the words that appear in both the changed regions of the diff and the target commit message, extracting their intersection. Subsequently, the generated message is decoded into a complete word sequence, and a zero mask matrix with the same expected length is initialized (Lines 4–5). Then, for each word in the generated message, the algorithm checks whether its lemmatized form exists in the intersection set. If the word appears in the intersection, the corresponding position in the mask is set to 1, indicating the parts of the generated message that are semantically aligned with the change regions (Lines 6–9). Through this process, the algorithm creates a mask (ca_mask) that highlights the alignment positions between the generated message and the changed code.

Next, we incorporate the change alignment mask into the loss function, guiding the model to focus on generating these important tokens correctly. The change alignment loss is computed as

$$\mathcal{L}_{CA} = - \sum_{i=1}^L ca_mask_i \log P(y_i | y_{<i}), \quad (7)$$

where L denotes the length of the generated sequence, and ca_mask_i is the value of the change alignment mask for the i -th token, corresponding to the probability of generating the i -th token in the sequence.

The final loss $\mathcal{L}_{\text{final}}$ is defined as the weighted sum of the change-aligned loss (\mathcal{L}_{CA}) and the standard cross-entropy loss (\mathcal{L}_{CE}):

$$\mathcal{L}_{\text{final}} = \lambda \cdot \mathcal{L}_{CA} + (1 - \lambda) \cdot \mathcal{L}_{CE}, \quad (8)$$

where λ is a hyperparameter controlling the balance between the two losses.

4 Evaluation

In this section, we conduct comprehensive experiments to evaluate the effectiveness of CAMCMG, aiming to address the following research questions.

- **RQ1: Overall performance.** How does CAMCMG perform compared to the state-of-the-art commit message generation methods?

Method	Java				C#				C++				Python				JavaScript								
	BLEU	Met.	Rou.	Cid.	BRSA	BLEU	Met.	Rou.	Cid.	BRSA	BLEU	Met.	Rou.	Cid.	BRSA	BLEU	Met.	Rou.	Cid.	BRSA					
NNGen (Liu et al., 2018)	19.41	12.40	25.15	1.23	29.40	20.93	13.52	25.14	1.43	29.42	12.72	8.48	16.85	0.67	21.71	15.62	10.42	20.97	0.88	25.46	17.96	11.88	23.47	1.15	27.78
RACE (Shi et al., 2022)	25.66	15.46	32.02	1.76	37.98	26.33	16.37	31.31	1.84	36.81	19.13	12.55	24.52	1.14	31.07	21.79	14.68	28.35	1.40	34.45	25.55	16.31	31.79	1.84	37.34
COME (He et al., 2023)	27.17	16.91	34.59	1.90	41.13	27.29	17.77	33.33	1.91	38.96	20.80	14.55	27.01	1.25	33.93	23.17	16.46	30.48	1.50	37.06	26.91	17.84	34.44	1.92	39.93
CodeT5-base (Wang et al., 2021b)	22.76	14.57	30.23	1.43	36.42	22.21	14.51	29.08	1.33	34.22	16.73	11.69	22.86	0.85	29.85	17.99	12.74	25.27	0.96	31.67	22.87	15.12	29.81	1.50	34.76
CCT5 (Lin et al., 2023)	26.46	16.58	33.95	1.81	40.81	<u>26.89</u>	<u>17.20</u>	<u>32.67</u>	<u>1.86</u>	<u>38.84</u>	19.80	13.82	26.11	1.15	33.37	<u>22.54</u>	<u>15.77</u>	<u>29.92</u>	<u>1.44</u>	<u>36.57</u>	26.16	17.17	33.44	1.86	39.12
Claude3 (The)	7.69	5.47	9.33	0.21	21.97	7.01	4.92	8.49	0.19	19.99	7.23	5.78	8.52	0.19	20.79	7.56	6.14	8.61	0.19	20.85	7.65	5.29	8.96	0.19	20.34
DeepSeekV2 (Liu et al., 2024)	6.56	4.88	7.96	0.21	21.47	5.94	4.25	6.96	0.18	19.03	6.47	5.13	7.27	0.19	19.95	6.84	5.45	7.41	0.20	20.05	6.41	4.65	7.24	0.19	19.34
Llama3 (Dubey et al., 2024)	7.88	6.29	9.63	0.25	24.13	7.25	5.69	8.58	0.21	21.58	7.97	7.08	8.95	0.24	22.86	8.43	7.43	9.15	0.25	22.71	7.64	5.97	8.99	0.23	21.78
CodeQwen1.5 (Bai et al., 2023)	7.37	6.49	8.93	0.21	22.54	7.09	5.81	7.99	0.17	20.51	7.56	6.74	8.34	0.19	21.38	7.93	7.01	8.18	0.20	21.11	6.87	6.08	8.01	0.18	20.06
GPT-4 (Achiam et al., 2023)	9.41	7.32	12.03	0.35	25.35	8.50	7.33	11.66	0.28	23.83	8.84	8.48	10.82	0.27	25.79	9.85	8.94	11.18	0.31	24.86	8.94	7.98	12.41	0.32	26.03
CAMCMG (Ours)	<u>27.10</u>	<u>17.63</u>	<u>38.30</u>	<u>2.06</u>	<u>41.90</u>	25.38	17.17	35.70	1.88	39.31	20.91	14.83	31.50	1.57	35.07	21.31	16.87	35.54	1.59	37.93	29.02	20.68	41.13	2.28	43.84

Table 1: The overall performance of CAMCMG. The best score is highlighted in **bold** and the second best score is underlined.

- RQ2: **Ablation study**. How does each component of CAMCMG contribute to the effectiveness?
- RQ3: **Human evaluation**. How does CAMCMG perform from the perspective of developers?

4.1 Evaluation Criteria

We evaluate the quality of the generated messages using five metrics: BLEU (Papineni et al., 2002) (the BLEU results reported in this paper are B-Norm), ROUGE (Lin, 2004), METEOR (Banerjee and Lavie, 2005), CIDEr (Vedantam et al., 2015) and BRSA (Li et al., 2024). The detailed metrics calculation can be found in Appendix A.

4.2 Dataset

The MCMD dataset, collected by Tao et al. (Tao et al., 2022), includes five programming languages: Java, C++, C#, Python, and JavaScript, and has been widely used in recent years for commit message generation (He et al., 2023; Shi et al., 2022; Tao et al., 2024). In this paper, we use the same version of the dataset as that used in COME (He et al., 2023) for our experiments.

4.3 Evaluation Models

We compare CAMCMG with one retrieval method (NNGen (Liu et al., 2018)), two learning-based methods (CodeT5 (Wang et al., 2021b), CCT5 (Lin et al., 2023)), two hybrid methods (COME (He et al., 2023), RACE (Shi et al., 2022)) and five LLM-based methods (Claude3-Haiku (The), DeepSeek-V2-16B (Liu et al., 2024), Llama3.1-8B (Dubey et al., 2024), CodeQwen1.5-7B (Bai et al., 2023), GPT-4 (Achiam et al., 2023)).

4.4 Experimental Setting

For the dataset, we set the maximum sequence length of *diff* to 512 tokens and commit messages

to 50 tokens. We initialize our model using the pre-trained CodeT5-base while training the encoder-decoder neural network. The original CodeT5 vocabulary contains 32,100 tokens, and following Shi et al. (Shi et al., 2022), we extend it with nine special tokens (<INSERT>, <INSERT_OLD>, <DELETE>, <DELETE_END>, <KEEP>, <KEEP_END>, <REPLACE_OLD>, <REPLACE_OLD>, <REPLACE_END>), resulting in a vocabulary size of 32,109. We use the AdamW optimizer with a learning rate of 5e-5, and the batch size is 16. The maximum number of epochs is set to 10. All experiments are conducted on an NVIDIA GeForce RTX 4090 GPU. For GPT-4, following prior works (Lopes et al., 2024; Wu et al., 2025), we randomly select 200 samples from the test set to reduce the number of API calls.

5 Results And Analysis

In this section, we first present the overall performance of CAMCMG (RQ1) in Section 5.1, followed by the results of the ablation study (RQ2) in Section 5.2, and then discuss the results of human evaluation (RQ3) in Section 5.3.

5.1 RQ1: Overall Performance

Table 1 presents the experimental results of different models on the MCMD dataset. As shown in Table 1, CAMCMG achieves the best performance, with an average improvement of 6.4% over the baselines. Specifically, METEOR improves by 4.2% on average, ROUGE-L by 14.1%, CIDEr by 11.4%, and BRSA by 3.7%. However, the BLEU score decreases by 1.4% on average. This is because BLEU relies on n-gram precision, so outputs that are semantically correct but lexically different may receive low scores (Wieting et al., 2019). In contrast, BRSA which measures semantic similarity remains higher, showing that the generated

Method	Java					C#					C++					Python					JavaScript				
	BLEU	Met.	Rou.	Cid.	BRSA	BLEU	Met.	Rou.	Cid.	BRSA	BLEU	Met.	Rou.	Cid.	BRSA	BLEU	Met.	Rou.	Cid.	BRSA	BLEU	Met.	Rou.	Cid.	BRSA
- LA	24.72	17.10	36.94	1.92	41.04	23.75	16.25	33.87	1.62	37.56	19.07	14.63	29.78	1.32	33.46	19.92	15.96	33.70	1.52	37.37	28.60	20.02	40.54	2.17	42.91
- CA	24.06	15.88	34.44	1.73	37.72	22.57	15.02	32.17	1.57	35.41	19.44	13.42	28.46	1.38	31.90	19.96	14.54	31.51	1.43	34.92	27.93	19.23	39.15	2.13	42.24
- CAL	26.04	16.44	36.39	1.89	39.81	22.68	15.12	32.19	1.58	35.70	19.79	13.96	29.26	1.41	32.77	20.98	15.34	32.91	1.50	34.80	28.03	19.38	38.96	2.16	42.20

Table 2: Results of the ablation study across different programming languages. LA, CA, and CAL denote Local Attention, Change Attention, and Change Alignment Loss, respectively. Like “-LA” indicate that the corresponding component has been removed in the ablation study.

messages are still semantically accurate despite the low BLEU.

It is worth noting that although LLM-based methods achieve relatively lower scores on automatic evaluation metrics, this does not necessarily indicate poor generation quality. Automatic evaluation metrics can measure the textual differences between generated messages and reference messages, but they are limited in capturing semantic consistency. Commit messages often have multiple expressions, and a generated message may use different words from the reference while conveying the same meaning, which can lead to underestimation by automatic metrics (Wang et al., 2024a). Therefore, to more accurately assess the quality of generated commit messages, we conducted a human evaluation in Section 5.3.

5.2 RQ2: Ablation Study

We conducted additional experiments to analyze the impact of different components in CAMCMG on the overall performance. Table 2 presents the performance comparison between the complete CAMCMG model and its ablated variants across different programming languages.

First, to investigate the impact of different attention mechanisms on the overall performance, we remove the local attention and the change attention separately. As shown in Table 2, when local and change attention are removed, the model’s average performance drops by approximately 5.3% and 9.6%, respectively. This indicates that both attention mechanisms positively contribute to the model’s performance, and the change attention has a more pronounced impact.

We further conduct experiments to explore the impact of the change alignment loss on the overall performance of CAMCMG. As shown in the last row of Table 2, removing this component leads to an average performance drop of approximately 7.1%. Specifically, BLEU decreases by 5.0%, METEOR by 8.0%, ROUGE-L by 6.9%, CIDEr by 9.1%, and BRSA by 6.5%. It indicates that the

change alignment loss plays a positive role in improving model performance.

In summary, the results clearly show that removing any of these components leads to a performance drop, indicating that each part contributes meaningfully to the model’s overall performance. In particular, the change attention mechanism yields the most substantial improvement, highlighting its importance in capturing cross-hunk semantic associations critical for commit message generation.

5.3 RQ3: Human Evaluation

Method	Informativeness	Conciseness	Expressiveness
NNGen	2.72	2.26	2.68
RACE	2.85	2.46	2.76
COME	2.87	2.59	2.83
CodeT5-base	2.85	2.54	2.84
CCT5	2.92	2.67	2.82
Claude3	3.17	3.03	3.01
DeepSeekV2	3.02	3.14	3.07
Llama3	3.10	2.98	3.12
CodeQwen1.5	2.99	3.05	3.14
GPT-4	3.34	3.29	3.38
CAMCMG (Ours)	<u>3.28</u>	<u>3.24</u>	<u>3.19</u>

Table 3: The result of human evaluation.

Following previous work (Liu et al., 2020; Wang et al., 2021a), we randomly select 100 commits from the testing set and design a questionnaire for manual evaluation. Each question in the questionnaire consisted of a code change, the corresponding reference message, and messages generated by various baseline models (The example of questionnaire can be seen in Appendix C). The study invited 2 PhD students and 3 Master students from the computer laboratories to participate in the survey, all of whom have at least 3 years of programming experience. Participants were asked to score each generated message on a scale of 1-5 across three dimensions: **Informativeness** (the extent to which the message reflects important information about the *diff*), **Conciseness** (the brevity and focus of the message), and **Expressiveness** (the clarity and naturalness in conveying the intent of the *diff*). The

performance of each model was measured by the average scores across all generated messages.

Table 3 presents the results of human evaluation. The results show that CAMCMG achieves the second-best overall performance, ranking slightly below GPT-4, while outperforming the best non-GPT baseline by 5.2%. Although GPT-4 attains higher scores, it incurs substantially higher costs in practical applications. Overall, CAMCMG can generate high-quality commit messages while maintaining lower usage costs. In addition, we conducted Wilcoxon signed-rank tests (Wilcoxon et al., 1963) on the human evaluation results to examine the significance of the differences between CAMCMG and other baselines. The results indicate that the performance differences are statistically significant, with all p-values smaller than 0.05 at the 95% confidence level.

6 Case Study

To further demonstrate the performance of our model in real scenarios, we analyze two cases and compare the generated messages with baseline methods.

```

@@ -521,7 +521,7 @@ public String getCode() {
    final int getCodeAtLength() {
-    return content() == null ? getCode().length() : content().getCodeAtLength();
+    return content().getCodeAtLength();
    }
}
@@ -748,7 +748,7 @@ public final LineLocation createLineLocation(int lineNumber) {
    * different source types.
    */
-    Object getHasKey() {
+    return content() == null ? getName() : content().getHasKey();
+    return content().getHasKey();
    }
}
final TextMap getHashMap() {
}

Commit message:
Reference: Always delegate to content() as it has to be always non-null
NWGen: Source - newFromTest is supported
RACE: merge pull request in g / truffle from ~ michael . van . de . vanter . oracle . com
COMe: merge pull request in g / truffle from paw / fix - source - equals to master
CodeTS-base: Merge pull request in G / truffle from bugfix / SourceImpl to master
CCTS: each instances of source needs its content ( ) .
Claude3: Refactor Source class to delegate equals and hashCode to Content
DeepSeekV2: Updated getCode and getHasKey methods to use newContent method
Llama3: Refactor Source.java to use content().equals() and content().hashCode()
CodeQwen: Refactor and optimize Source class for better performance and readability.
GPT-4: Simplify Source methods by removing unnecessary null handling for content()
CAMCMG: remove redundant null checks for content() in Source.equals() and hashCode()

```

Figure 3: Example of removing redundant null checks.

Figure 3 presents a real-world case (the full diff can be found on Appendix E), which includes the code changes, the reference commit message, as well as the commit messages generated by CAMCMG and other baseline models. In this case, the developers refactored the Source class by removing redundant null checks related to the content() method. The manually written commit message accurately captures the developers' intent. The diff snippet shown in the figure contains multiple hunks across different functions, all aimed at removing unnecessary null checks for content(). From the generated results, only CAMCMG and Llama3 produce commit messages that are largely consis-

tent with the reference message. In contrast, the outputs from other baseline models show varying degrees of deviation from the expected results.

```

@@ -263,6 +263,9 @@ public Message deserialize(JsonElement json, Type typeOfT,
    Message.Builder protoBuilder =
    (Message.Builder) getCacheMethod(protoClass, "newBuilder").invoke(null);
+
+    Message defaultInstance =
+    (Message) getCacheMethod(protoClass, "getDefaultInstance").invoke(null);
+
    Descriptor protoDescriptor =
    (Descriptor) getCacheMethod(protoClass, "getDescriptor").invoke(null);
    // Call setters on all of the available fields
@@ -300,8 +303,7 @@ public Message deserialize(JsonElement json, Type typeOfT,
    fieldValue = context.deserialize(jsonElement, protoArrayFieldType);
    protoBuilder.setField(fieldDescriptor, fieldValue);
    } else {
-    Message prototype = protoBuilder.build();
-    Object field = prototype.getField(fieldDescriptor);
+    Object field = defaultInstance.getField(fieldDescriptor);
+    fieldValue = context.deserialize(jsonElement, field.getClass());
    protoBuilder.setField(fieldDescriptor, fieldValue);
    }
}

Commit message:
Reference: avoid repeated calls to build() during deserialization by obtaining field types
from the Message's defaultInstance
NWGen: APT Change : Implemented various getAs methods to return null .
RACE: allow subclasses to provide their own default instance .
COMe: use a cached default instance of the proto type adapter .
CodeTS-base: Use getDefaultInstance ( ) instead of newBuilder ( ) .
CCTS: fix deserialization of default instance messages .
Claude3: Add support for Message.setField with default instance and prototype
DeepSeekV2: Update ProtobufAdapter to use getDefaultInstance and setters for deserialization
Llama3: Improved ProtobufAdapter to handle null fields and default instances
CodeQwen: Implement a custom Gson type adapter for Protobuf messages.
GPT-4: Optimize deserialization by using defaultInstance instead of building a prototype
instance for field access.
CAMCMG: Use defaultInstance instead of building a new prototype in deserialization

```

Figure 4: Example of instance reuse in deserialization.

To further verify the model's capability in understanding semantic changes in code, we conducted a case study on a real-world commit related to "protobuf message deserialization optimization". As shown in Figure 4, this commit replaces the invocation of protoBuilder.build() with getDefaultInstance(), thereby avoiding repeated construction of temporary prototype objects during deserialization and improving execution efficiency. In this case, our model successfully links the two hunks, accurately identifying their semantic association and generating a commit message that better reflects the developer's actual intent. In contrast, baseline models only capture partial information, resulting in incomplete semantic understanding in their generated commit messages.

7 Conclusion

In this paper, we propose CAMCMG, a change-aware model for commit message generation that consists of three core components. First, a local attention mechanism is designed to capture local contextual information. Second, a change-aware attention mechanism is introduced to model semantic associations across different hunks. Third, a change alignment loss is employed to explicitly align the diff with its commit message. In addition, we conduct a human evaluation to assess the quality of the generated commit messages, demonstrating the effectiveness of CAMCMG. Our code and data are publicly available at <https://anonymous.4open.science/r/CAMCMG/>.

615 Limitations

616 We have identified the following main limitations:

617 *Generalizability to other programming lan-*
618 *guages:* Our experiments are conducted on a lim-
619 ited set of programming languages. The generaliz-
620 ability of the proposed approach to other program-
621 ming languages remains an open question.

622 *Metrics:* The model may generate commit mes-
623 sages that are semantically correct but lexically
624 diverse, resulting in relatively low scores on the
625 automatic metrics. Although BRSA can capture
626 semantic similarity, the limitations of automatic
627 metrics remain an issue for evaluation.

628 *Change Alignment Loss:* The change alignment
629 loss relies on the presence of tokens in the *diff* that
630 correspond to tokens in the commit message. If no
631 such tokens exists, its effectiveness may be limited.

632 References

633 [The claude 3 model family: Opus, sonnet, haiku.](#)

634 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
635 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
636 Diogo Almeida, Janko Altenschmidt, Sam Altman,
637 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-
638 cal report. *arXiv preprint arXiv:2303.08774*.

639 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
640 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
641 Huang, and 1 others. 2023. Qwen technical report.
642 *arXiv preprint arXiv:2309.16609*.

643 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An
644 automatic metric for mt evaluation with improved cor-
645 relation with human judgments. In *Proceedings of*
646 *the acl workshop on intrinsic and extrinsic evaluation*
647 *measures for machine translation and/or summariza-*
648 *tion*, pages 65–72.

649 Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020.
650 Longformer: The long-document transformer. *arXiv*
651 *preprint arXiv:2004.05150*.

652 Raymond PL Buse and Westley R Weimer. 2010. Au-
653 tomatically documenting program changes. In *Pro-*
654 *ceedings of the 25th IEEE/ACM international con-*
655 *ference on automated software engineering*, pages
656 33–42.

657 Maria Caulo, Bin Lin, Gabriele Bavota, Giuseppe Scan-
658 niello, and Michele Lanza. 2020. Knowledge transfer
659 in modern code review. In *Proceedings of the 28th In-*
660 *ternational Conference on Program Comprehension*,
661 pages 230–240.

662 Gonçalo M Correia, Vlad Niculae, and André FT Mar-
663 tins. 2019. Adaptively sparse transformers. *arXiv*
664 *preprint arXiv:1909.00015*.

665 Michael Denkowski and Alon Lavie. 2014. Meteor
666 universal: Language specific translation evaluation
667 for any target language. In *Proceedings of the ninth*
668 *workshop on statistical machine translation*, pages
669 376–380.

670 Jinhao Dong, Yiling Lou, Qihao Zhu, Zeyu Sun, Zhilin
671 Li, Wenjie Zhang, and Dan Hao. 2022. Fira: fine-
672 grained graph-based code change representation for
673 automated commit message generation. In *Proceed-*
674 *ings of the 44th international conference on software*
675 *engineering*, pages 970–981.

676 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
677 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
678 Akhil Mathur, Alan Schelten, Amy Yang, Angela
679 Fan, and 1 others. 2024. The llama 3 herd of models.
680 *arXiv e-prints*, pages arXiv–2407.

681 Robert Dyer, Hoan Anh Nguyen, Hridayesh Rajan, and
682 Tien N Nguyen. 2013. Boa: A language and in-
683 frastructure for analyzing ultra-large-scale software
684 repositories. In *2013 35th International Conference*
685 *on Software Engineering*, pages 422–431. IEEE.

686 Tudor Gîrba and Stéphane Ducasse. 2006. Modeling
687 history to analyze software evolution. *Journal of*
688 *Software Maintenance and Evolution: Research and*
689 *Practice*, 18(3):207–236.

690 Daya Guo, Canwen Xu, Nan Duan, Jian Yin, and Ju-
691 lian McAuley. 2023. Longcoder: A long-range pre-
692 trained language model for code completion. In *In-*
693 *ternational Conference on Machine Learning*, pages
694 12098–12107. PMLR.

695 Yichen He, Liran Wang, Kaiyi Wang, Yupeng Zhang,
696 Hang Zhang, and Zhoujun Li. 2023. Come: Commit
697 message generation with modification embedding.
698 In *Proceedings of the 32nd ACM SIGSOFT Interna-*
699 *tional Symposium on Software Testing and Analysis*,
700 pages 792–803.

701 Thong Hoang, Hong Jin Kang, David Lo, and Julia
702 Lawall. 2020. Cc2vec: Distributed representations of
703 code changes. In *Proceedings of the ACM/IEEE 42nd*
704 *international conference on software engineering*,
705 pages 518–529.

706 Yuan Huang, Nan Jia, Hao-Jie Zhou, Xiang-Ping Chen,
707 Zi-Bin Zheng, and Ming-Dong Tang. 2020. Learning
708 human-written commit messages to document code
709 changes. *Journal of Computer Science and Technol-*
710 *ogy*, 35(6):1258–1277.

711 Aaron Imani, Iftexhar Ahmed, and Mohammad
712 Moshirpour. 2024. Context conquers parameters:
713 Outperforming proprietary llm in commit message
714 generation. *arXiv preprint arXiv:2408.02502*.

715 Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas
716 Pfister. 2021. Cutpaste: Self-supervised learning for
717 anomaly detection and localization. In *Proceedings*
718 *of the IEEE/CVF conference on computer vision and*
719 *pattern recognition*, pages 9664–9674.

720	Cong Li, Zhaogui Xu, Peng Di, Dongxia Wang, Zheng Li, and Qian Zheng. 2024. Understanding code changes practically with small-scale language models. In <i>Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering</i> , pages 216–228.	773
721		774
722		775
723		776
724		777
725		
726	Yingling Li, Yuhan Wu, Zi'ao Wang, Lei Huang, Junjie Wang, Jianping Li, and Mingying Huang. 2025. Code-doctor: multi-category code review comment generation. <i>Automated Software Engineering</i> , 32(1):25.	778
727		779
728		780
729		781
730	Bo Lin, Shangwen Wang, Zhongxin Liu, Yepang Liu, Xin Xia, and Xiaoguang Mao. 2023. Cct5: A code-change-oriented pre-trained model. In <i>Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering</i> , pages 1509–1521.	782
731		783
732		784
733		785
734		786
735		787
736	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	788
737		789
738		790
739	Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, and 1 others. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. <i>arXiv preprint arXiv:2405.04434</i> .	791
740		792
741		793
742		794
743		795
744		
745	Shangqing Liu, Cuiyun Gao, Sen Chen, Lun Yiu Nie, and Yang Liu. 2020. Atom: Commit message generation based on abstract syntax tree and hybrid ranking. <i>IEEE Transactions on Software Engineering</i> , 48(5):1800–1817.	796
746		797
747		798
748		799
749		800
750	Zhongxin Liu, Zhijie Tang, Xin Xia, and Xiaohu Yang. 2023. Ccrep: Learning code change representations via pre-trained code model and query back. In <i>2023 IEEE/ACM 45th International Conference on Software Engineering</i> , pages 17–29. IEEE.	801
751		802
752		803
753		804
754		805
755	Zhongxin Liu, Xin Xia, Ahmed E Hassan, David Lo, Zhenchang Xing, and Xinyu Wang. 2018. Neural-machine-translation-based commit message generation: how far are we? In <i>Proceedings of the 33rd ACM/IEEE international conference on automated software engineering</i> , pages 373–384.	806
756		807
757		808
758		809
759		810
760		811
761	Cristina V Lopes, Vanessa I Klotzman, Iris Ma, and Iftekar Ahmed. 2024. Commit messages in the age of large language models. <i>arXiv preprint arXiv:2401.17622</i> .	812
762		813
763		814
764		815
765	Walid Maalej and Hans-Jörg Happel. 2010. Can development work describe itself? In <i>2010 7th IEEE working conference on mining software repositories</i> , pages 191–200. IEEE.	816
766		817
767		818
768		819
769	Lun Yiu Nie, Cuiyun Gao, Zhicong Zhong, Wai Lam, Yang Liu, and Zenglin Xu. 2021. Coregen: Contextualized code representation learning for commit message generation. <i>Neurocomputing</i> , 459:97–107.	820
770		821
771		822
772		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900

827 Chuangwei Wang, Li Zhang, and Xiaofang Zhang. 884
828 2024a. Multi-grained contextual code representation 885
829 learning for commit message generation. *Informa- 886*
830 *tion and Software Technology*, 167:107393.

831 Haoye Wang, Xin Xia, David Lo, Qiang He, Xinyu 887
832 Wang, and John Grundy. 2021a. Context-aware 888
833 retrieval-based deep commit message generation. 889
834 *ACM Transactions on Software Engineering and 890*
835 *Methodology (TOSEM)*, 30(4):1–30. 891

836 Hongqi Wang, Jiawen Zhao, Qunpo Liu, Naohiko Hana-
837 jima, and Xuhui Bu. 2025. Gfd-net: Micro-electrical
838 connector defect localization network with global addi-
839 tive attention and frequency-aware feature fusion.
840 *Signal, Image and Video Processing*, 19(11):890.

841 Yanlin Wang, Yanxian Huang, Daya Guo, Hongyu
842 Zhang, and Zibin Zheng. 2024b. Sparsecoder:
843 Identifier-aware sparse transformer for file-level code
844 summarization. In *2024 IEEE International Confer-
845 ence on Software Analysis, Evolution and Reengi-
846 neering*, pages 614–625. IEEE.

847 Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH
848 Hoi. 2021b. Codet5: Identifier-aware unified
849 pre-trained encoder-decoder models for code un-
850 derstanding and generation. *arXiv preprint*
851 *arXiv:2109.00859*.

852 Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin
853 Carbonneaux, Lingming Zhang, Daniel Fried,
854 Gabriel Synnaeve, Rishabh Singh, and Sida I Wang.
855 2025. Swe-rl: Advancing llm reasoning via reinforce-
856 ment learning on open software evolution. *arXiv*
857 *preprint arXiv:2502.18449*.

858 John Wieting, Taylor Berg-Kirkpatrick, Kevin Gimpel,
859 and Graham Neubig. 2019. Beyond bleu: training
860 neural machine translation with semantic similarity.
861 *arXiv preprint arXiv:1909.06694*.

862 Frank Wilcoxon, SK Katti, Roberta A Wilcox, and 1
863 others. 1963. *Critical values and probability levels*
864 *for the Wilcoxon rank sum test and the Wilcoxon*
865 *signed rank test*, volume 1. American Cyanamid
866 Pearl River, NY.

867 Yifan Wu, Yunpeng Wang, Ying Li, Wei Tao, Siyu
868 Yu, Haowen Yang, Wei Jiang, and Jianguo Li. 2025.
869 An empirical study on commit message generation
870 using llms via in-context learning. *arXiv preprint*
871 *arXiv:2502.18904*.

872 Manzil Zaheer, Guru Guruganesh, Kumar Avinava
873 Dubey, Joshua Ainslie, Chris Alberti, Santiago On-
874 tanon, Philip Pham, Anirudh Ravula, Qifan Wang,
875 Li Yang, and 1 others. 2020. Big bird: Transformers
876 for longer sequences. *Advances in neural informa-
877 tion processing systems*, 33:17283–17297.

878 Jian Zhang, Chong Wang, Anran Li, Wenhan Wang,
879 Tianlin Li, and Yang Liu. 2024a. Vuladvisor: Natural
880 language suggestion generation for software vulnera-
881 bility repair. In *Proceedings of the 39th IEEE/ACM*
882 *International Conference on Automated Software En-
883 gineering*, pages 1932–1944.

Linghao Zhang, Hongyi Zhang, Chong Wang, and Peng
Liang. 2024b. Rag-enhanced commit message gener-
ation. *arXiv preprint arXiv:2406.05514*.

Yuxia Zhang, Zhiqing Qiu, Klaas-Jan Stol, Wenhui Zhu,
Jiaxin Zhu, Yingchen Tian, and Hui Liu. 2024c. Au-
tomatic commit message generation: A critical re-
view and directions for future work. *IEEE Transac-
tions on Software Engineering*, 50(4):816–835.

A Evaluation Metrics

BLEU (Papineni et al., 2002): BLEU is an evaluation metric based on n-gram precision. It calculates a score by comparing the overlap of n-grams between the generated text and the reference text. Traditional BLEU performs poorly on short texts. Therefore, in this paper, we adopt a variant of BLEU, B-NORM (Tao et al., 2021) (the BLEU results reported in this paper are B-Norm), which introduces a normalization strategy suitable for short texts such as commit messages, allowing a more reasonable reflection of the model’s actual performance.

$$BLEU - N = BP \cdot \exp\left(\sum_{n=1}^N w_n \cdot \log(p_n)\right),$$

where BP denotes the length penalty, which penalizes predictions that are too short. p_n represents the empirically chosen weight for each n -gram, corresponding to the matching score of the n -grams. The score is calculated as follows:

$$BP = \begin{cases} 1 & \text{if } c > 0 \\ e^{(1-r/c)} & \text{if } c \leq 0 \end{cases},$$

where c and r denote the lengths of the generated text and the reference text, respectively.

ROUGE (Lin, 2004): A metric focuses on recall rather than precision. It measures how much information in the reference text is covered by the generated text. ROUGE-L calculates the score based on the precision and recall of the Longest Common Subsequence (LCS) between the generated and reference texts. The score is calculated as follows:

$$ROUGE-L = \frac{(1 + \beta^2)R_{lcs}P_{lcs}}{R_{lcs} + \beta^2P_{lcs}}$$

where R_{lcs} and P_{lcs} represent the recall and precision, respectively. In DUC, β is set to a very large value, so only R_{lcs} is effectively considered.

The recall and precision are defined as:

$$R_{lcs} = \frac{LCS(X, Y)}{m}, \quad P_{lcs} = \frac{LCS(X, Y)}{n}$$

where $LCS(X, Y)$ denotes the length of the longest common subsequence between X and Y , and m and n are the lengths of the reference text and the generated text, respectively.

METEOR (Banerjee and Lavie, 2005): A word-based metric used to measure the extent to which

the generated text covers the reference text. It calculates a score by matching words between the generated and reference texts. Higher scores indicate greater similarity between the generated and reference texts, and hence better quality. The score is calculated as follows:

$$METEOR = (1 - Pen)F_{mean}$$

where Pen is calculated according to the number of chunks (ch) and the number of matches (m):

$$Pen = \gamma * \left(\frac{ch}{m}\right)^\beta$$

The values of β and γ parameters are set to 0.20 and 0.60, respectively, following the prior work (Denkowski and Lavie, 2014).

CIDEr (Vedantam et al., 2015): It measures the similarity between the generated text and reference texts by computing TF-IDF weighted n-gram similarities, and then averages the scores across different n-grams to obtain the final evaluation result. The score is calculated as follows:

$$CIDEr_n(c_i, s_i) = \frac{\langle g^n(c_i), g^n(s_i) \rangle}{\|g^n(c_i)\| \|g^n(s_i)\|},$$

where $g^n(s_i)$ is the vector formed by $g_k(s_i)$, containing all n-grams (n ranges from 1 to 4). c_i denotes the i -th generated sentence. Finally, the overall CIDEr score is computed by combining the scores from different n-grams, as follows:

$$CIDEr(c_i, s_i) = \sum_{n=1}^N \omega_n CIDEr_n(c_i, s_i)$$

BRSA (Li et al., 2024): This metric utilizes MP-Net (Song et al., 2020), specifically *all-mpnet-base-v2* to compute the semantic similarity between the reference and generated texts, denoted as SEMSIM. To mitigate the limitations imposed by BLEU and ROUGE-L, Song et al. (Song et al., 2020) apply a double weighting on SEMSIM and combine it with ROUGE-L and BLEU using weights of 0.25, 0.25, and 0.5, respectively, to calculate the BRSA score.

$$BRSA = 0.25 \times BLEU + 0.25 \times ROUGE-L + 0.5 \times SEMSIM$$

B Dataset

Table 4 presents the statistics of commits and their corresponding hunks across five programming languages and three dataset splits (Training, Valid,

943
944
945

and Test). It can be observed that, on average, each commit contains approximately 3 to 5 hunks. These statistics provide a foundational basis for the subsequent analysis and modeling in this study.

Language	Phase	Commits	Hunks	Avg
Java	Training	160,018	591,633	3.70
	Valid	19,825	74,294	3.75
	Test	20,159	78,714	3.90
C#	Training	149,907	677,062	4.52
	Valid	18,688	87,777	4.70
	Test	18,702	81,009	4.33
C++	Training	160,948	739,925	4.60
	Valid	20,000	95,139	4.76
	Test	20,141	87,922	4.37
Python	Training	206,777	953,388	4.61
	Valid	25,912	128,373	4.95
	Test	25,837	117,231	4.54
JavaScript	Training	197,529	806,673	4.08
	Valid	24,899	120,697	4.85
	Test	24,773	92,523	3.73

Table 4: Statistics of the dataset.

946

C Questionnaire of Human Evaluation

We designed a survey for human evaluation in which each participant was asked to score the commit messages generated by each model with respect to a given code change across three dimensions: informativeness, conciseness, and expressiveness. Scores ranged from 1 to 5, with higher values indicating better quality. Figure 5 illustrates an example question from the survey, where the order of the commit messages generated by different models was randomized to avoid bias.

947

948

949

950

951

952

953

954

955

956

```

Raw Diff
@@ -521,7 +521,7 @@ public String getCode() {
     final int getCodeLength() {
-        return content() == null ? getCode().length() : content().getCodeLength();
+        return content().getCodeLength();
     }
@@ -748,7 +748,7 @@ public final LineLocation createLineLocation(int lineNumber) {
     * different source types.
     }
     Object getHashKey() {
-        return content() == null ? getName() : content().getHashKey();
+        return content().getHashKey();
     }
     final TextMap getTextMap() {
Commit message:
Reference: Always delegate to content() as it has to be always non-null
Generate Message1: Source . newFromTest is supported
informativeness 1 2 3 4 5 conciseness 1 2 3 4 5 expressiveness 1 2 3 4 5
Generate Message2: remove redundant null checks for content() in Source.equals() and hashCode()
informativeness 1 2 3 4 5 conciseness 1 2 3 4 5 expressiveness 1 2 3 4 5
Generate Message3: merge pull request in g / truffile from paw / fix - source - equals to master
informativeness 1 2 3 4 5 conciseness 1 2 3 4 5 expressiveness 1 2 3 4 5
Generate Message4: Merge pull request in G / truffile from bugfix / SourceImpl to master
informativeness 1 2 3 4 5 conciseness 1 2 3 4 5 expressiveness 1 2 3 4 5
Generate Message5: each instances of source needs its content() .
informativeness 1 2 3 4 5 conciseness 1 2 3 4 5 expressiveness 1 2 3 4 5

```

Figure 5: The example of questionnaire.

957

D Hyperparameter Experiments

958

In addition, we further conducted experiments on two key hyperparameters, namely the window size and the weighting factor λ , to investigate their effects on model performance. For the window size parameter, we set the values to 32, 64, 128, and 256 to analyze how different local context ranges influence the effectiveness of the sparse attention mechanism. As shown in Figure 6(a), the model achieves its best performance when the window size is set to 64. When the window size is too small, the model lacks sufficient contextual information to understand a code change. Conversely, an excessively large window introduces redundant context and increases attention noise, thereby diluting the model’s focus on essential regions.

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

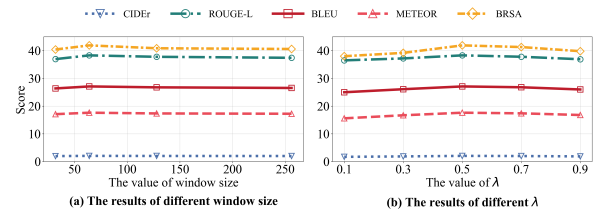


Figure 6: The results of hyperparameter experiments.

For the parameter λ , we vary its value from 0.1 to 0.9 with a step size of 0.2 to examine its role in balancing the change alignment loss and the standard cross-entropy loss (Equation 8). A smaller λ causes the model to rely more on the standard cross-entropy loss, which tends to optimize the overall linguistic quality of the generated messages, while a larger λ emphasizes the change alignment loss, thus enhancing the semantic alignment between the *diff* and its commit messages. As shown in Figure 6(b), the model performs best at $\lambda = 0.5$, suggesting that assigning a moderate weight to the change alignment loss can effectively improve the semantic consistency between code changes and commit messages.

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

Answering to hyperparameter experiments:
We conducted hyperparameter experiments on the window size and λ , and the results show that the model achieves its best performance when $\lambda = 0.5$ and the window size is set to 64.

990

E Full Case

991

In Section 6, due to the excessive length of the *diff*, only a partial *diff* is shown in the case1. The complete *diff* can be found in the Figure 7.

992

993

994

```

@@ -521,7 +521,7 @@ public String getCode() {
-   final int getCodeLength() {
+   return content() == null ? getCode().length() : content().getCodeLength();
+   return content().getCodeLength();
-   }
-748,7 +748,7 @@ public final LineLocation createLineLocation(int lineNumber) {
+   /*
+   Object getHasKey() {
-   return content() == null ? getName() : content().getHasKey();
+   return content().getHasKey();
-   }
-   final TextMap getTextMap() {
@@ -815,29 +815,23 @@ final boolean equalMime(Source other) {
-   return mimeType.equals(other.mimeType);
-   }
-   @Override
-   public boolean equals(Object obj) {
-   if (obj instanceof Source) {
-   Source other = (Source) obj;
-   if (content() == null) {
-   return super.equals(obj);
-   }
-   return content().equals(other.content()) && equalMime(other);
+   private static class Impl extends Source implements Cloneable {
+   Impl(Content content) {
+   super(content);
-   }
-   return false;
-   }
-   @Override
-   if (content() == null) {
-   return super.hashCode();
+   @Override
+   public boolean equals(Object obj) {
+   if (obj instanceof Source) {
+   Source other = (Source) obj;
+   return content().equals(other.content()) && equalMime(other);
+   }
+   return false;
-   }
-   return content().hashCode();
-   }
-   private static class Impl extends Source implements Cloneable {
-   Impl(Content content) {
-   super(content);
+   @Override
+   public int hashCode() {
+   return content().hashCode();
-   }
-   }
}

```

Commit message:

Reference: Always delegate to content() as it has to be always non-null
 NNGen: Source.newFromTest is supported
 RACE: merge pull request in g / truffle from ~ michael . van . de . vanter_oracle . com
 COME: merge pull request in g / truffle from paw / fix - source - equals to master
 Code15-base: Merge pull request in g / truffle from bugfix / SourceImpl to master
 CCTS: each instances of source needs its content ()
 Claude3: Refactor Source class to delegate equals and hashCode to Content
 DeepSeekV2: Updated getCode and getHasKey methods to use newContent method
 llama3: Refactor Source.java to use content().equals() and content().hashCode()
 CodeQwen: Refactor and optimize Source class for better performance and readability.
 GPT-4: Simplify Source methods by removing unnecessary null handling for content()
 CAMCMG: remove redundant null checks for content() in Source.equals() and hashCode()

Figure 7: The complete *diff* of case 1.