TOO NOISY TO LEARN: ENHANCING DATA QUALITY FOR CODE REVIEW COMMENT GENERATION

10/8/25

Chunhua Liu, Hong Yi Lin, Patanamon Thongtanunam,

The University of Melbourne

Xiangrui Ke,

David R. Cheriton School of Computer Science



EXSII

Existin

or no

```
@@ -80,6 +80,7 @@ public class HoodieCreateHandle<T extends
HoodieRecordPayload> extends HoodieIOH
String partitionPath, String fileId, Iterator<HoodieRecord<T>>
    recordIterator) {
    this(config, commitTime, hoodieTable, partitionPath, fileId);
    this.recordIterator = recordIterator;
    this.useWriterSchema = true;
```

 Model Reviewer's Comment: Why do we have this flag? comm Label: Noisy Comment

```
    Review
```

Reviewer's Comment: This can be simplified as new ArrayList<>(

Arrays.asList(new ProtocolConfig(protocol)))

Label: Valid Comment

WATERLOO

BACKGROUND

- Code review is essential but time-consuming and inconsistent.
- AI models can help generate review comments.
- Training datasets are **noisy** (vague, difficult-to-understand)
 → poor model outputs.

PRESENTATION TITLE

PAGE 2



IDEA

PRESENTATION TITLE

Use **Large Language Models (LLMs)** to semantically detect and filter out noisy comments.



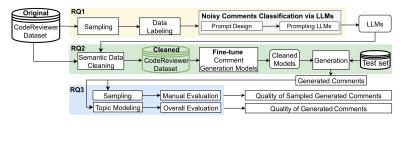
- Manually label a sample (valid or noisy?)
- Use LLMs for classification
- Train review generation models with high-quality cleaned data







STUDY DESIGN



RQ.1

_

RQ.3

To what degree can large language models semantically clean code review comments?

Does semantic data cleaning impact the accuracy of code review comment generation models?

RQ.2

Does semantic data cleaning improve the quality of code review comment generation models?

PRESENTATION TITLE PAGE 5



RQ2: Clean Data Impact the Accuracy of Generation

TABLE III: Model Performance (BLEU-4) on Comment Generation Models.

M	Dataset	Test	Valid _{Our&Tufano}	Noisy _{Our&Tufano}	Valid _{Our}	Noisy _{Our}	Valid _{Tufano}	Noisy _{Tufano}
CodeReviewer	ORIGINAL CLEANEDGPT-3.5 CONTROLLEDGPT-3.5 CLEANED _{LLAMA3} CONTROLLED _{LLAMA3}	5.73 6.04 5.4%↑* 5.63 5.97 4.2%↑*	6.17 6.97 13.0 %↑* 6.20 <u>6.63</u> 7.5%↑*	5.41 5.02 7.2%↓ 5.43 5.18 4.3%↓ 5.66	5.45 5.93 8.8%↑ 5.21 5.64 3.5%↑ 5.12	5.17 5.19 0.4%↑ 5.13 5.11 1.2%↓ 5.36	7.12 7.99 12.2%↑* 7.39 7.71 8.3%↑* 7.45	5.60 4.83 13.8%↓ 5.70 5.14 8.2%↓ 5.86
CodeT5	ORIGINAL CLEANEDGPT-3.5 CONTROLLEDGPT-3.5 CLEANEDLLAMA3 CONTROLLEDLLAMA3	5.19 5.67 9.2 %†* 5.20 5.54 <u>6.7%</u> †* 5.21	5.34 6.00 <u>12.4%</u> ↑* 5.34 5.74 7.5%↑* 5.19	5.04 5.23 3.8%↑* 5.30 5.33 5.8%↑ 5.12	4.84 <u>5.88</u> 21.5 %↑* 5.17 5.32 <u>9.9%</u> ↑* 4.95	5.09 5.27 3.5%↑ 5.39 5.14 1.0%↑ 5.26	5.85 6.06 3.6%↑ 5.45 6.09 4.1%↑ 5.38	6.03 5.15 14.6%↓ 5.41 5.46 9.5%↓ 5.01

The highest and second-highest results are in bold and underlined, respectively. * indicates the statistical significance (p-value < 0.05).

RQ1: LLMS Semantically Clean Code Review Comments

Prompt	Model	Input	Overall (weighted)		Valid (172)			Noisy (98)					
			Prec	Rec	F1	Prec	Rec	F1	#	Prec	Rec	F1	#
	Baseline [5]	-	40.6	63.7	49.6	63.7	100	77.8	270	0	0	0	0
ITION	GPT-3.5 CodeLlama Llama3	$R_{ m NL}$ $R_{ m NL}$ $R_{ m NL}$	70.3 64.1 71.8	54.1 65.6 72.6	55.7 58.0 71.7	85.1 66.0 75.3	36.6 94.8 84.9	51.2 77.8 79.8	74 247 194	44.4 60.9 65.8	88.8 14.3 51	59.2 23.1 57.5	196 23 76
PDEFINITION	GPT-3.5 CodeLlama Llama3	$R_{ m NL}$ + $C_{ m DIFF}$ $R_{ m NL}$ + $C_{ m DIFF}$ $R_{ m NL}$ + $C_{ m DIFF}$	65.6 54.2 62.6	61.5 62.2 65.2	62.2 52.3 59.8	75.8 63.8 66.7	58.1 94.2 90.7	65.8 76.1 76.8	132 254 234	47.8 37.5 55.6	67.3 6.1 20.4	55.9 10.5 29.9	138 16 36
JARY	GPT-3.5 CodeLlama Llama3	$R_{ m NL} \ R_{ m NL} \ R_{ m NL}$	66.8 71.0 71.0	59.2 71.7 71.9	59.5 70.1 <u>70.6</u>	49.7 73.2 74.0	60.6 87.7 86.0	54.6 79.8 79.6	107 205 200	46.2 67.2 65.7	76.3 43.9 46.9	57.6 53.1 54.8	160 64 70
Раскішаку	GPT-3.5 CodeLlama Llama3	$R_{ m NL}$ + $C_{ m DIFF}$ $R_{ m NL}$ + $C_{ m DIFF}$ $R_{ m NL}$ + $C_{ m DIFF}$	60.4 47.0 63.3	55.9 62.2 65.6	56.7 55.7 62.2	70.5 64.1 68.0	52.9 92.4 86.6	60.5 75.7 76.2	129 248 219	42.6 40.9 54.9	61.2 9.2 28.6	50.2 15.0 37.6	141 22 51

The highest and second-highest results are in bold and underlined. # represents the number of instances predicted in each class.

LLMs is a good classifier.

PAGE 6

PRESENTATION TITLE



RQ3: Clean Data Improve the Quality of Generation

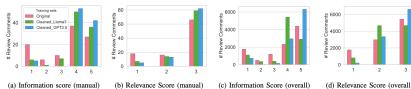


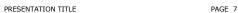
Fig. 4: Distribution of information and relevance scores on tests across CodeReviewer models trained on different training sets.

The cleaned models can generate more informative and relevant comments



Fig. 5: Example comments generated by original and cleaned models with information (Info) and Relevance (Rel) scores.





CONCLUSION

- To RQ1: LLM-based approach achieves 66-85% precision in detecting valid comments
- To RQ2: Using the predicted valid comments to fine-tune the state-of-the-art code review models (cleaned models)
- To RQ3: Cleaned models can generate more informative and relevant comments than the original models

Quality > Quantity

PRESENTATION TITLE

PAGE 9



PROS

- First to use LLMs for semantic cleaning of code review datasets
- Comprehensive experiment with different models and metrics.
- Insightful Resulte: better data quality improves model performance.
- Pioneeringly investigate the feasibility of using LLMs.

PRESENTATION TITLE

PAGE 10



CONS

- Bias exists
- Manual job still require
- Limited generalizability

RATING 4/5

- LLM-based cleaning significantly improve model quality
- Come out with the key insight: Quality > Quantity
- Advanced LLM-assisted exploration in software engineering





FUTURE WORK

- Expand to both validity and correctness.
- Bring more LLM-driven work to reduce manual effort
- More dataset and ablation experiments
- Domain generalization

PRESENTATION TITLE

PAGE 13



PRESENTATION TITLE

PAGE 14

DISCUSSION

10/8/





Thank you

PRESENTATION TITLE PAGE 15

10/8/