Impact of Evaluation Methodologies on Code Summarization

Pengyu Nie, Jiyang Zhang, Junyi Jessy Li, Raymond J. Mooney, Milos Gligoric

The University of Texas at Austin



ACL 2022

partially supported by Google





ML Models for Code Summarization

 Advances in ML/NLP are helping developers to produce and maintain software-related artifacts, for example code summarization

```
public static String selectText(XPathExpression expr, Node context) {
 try {
   return (String)expr.evaluate(context, XPathConstants.STRING );
 } catch (XPathExpressionException e) {
                                                                                 /** Evaluates the xpath expression as text. */
   throw new XmlException(e);
                                                                comment generation
               int ?(String str, char ch) {
                 int num = 0;
                 int index = -1;
                 do {
                   index = str.index0f(ch, index + 1);
                                                                                  count0ccurrences
                   if (index \geq 0) num++;
                                                                   method naming
                 } while (index >= 0);
                 return num;
```

Evaluation Methodologies of Code Summarization ML Models

- extract a dataset of (code, comment) samples
- split the dataset into training, validation, test sets
- train on training + validation sets
- report automatic metrics on test set





Temporal Relations Not Explicitly Modeled in Prior Work

written in 2018

```
/** Returns total number of connections in the pool. */
public synchronized int connectionCount() {
   return connections.size();
}
```



evaluation methodologies in prior work

2 written in 2019

/** Returns the number of idle connections in the pool. */
public synchronized int idleConnectionCount() {
 int total = 0;
 for (RealConnection connection : connections) {
 if (connection.allocations.isEmpty()) total++;
 }
 return total;
}





misunderstanding if a model might be useful / not useful once adopted



Our Contributions

- Study the evaluation methodologies of 18 recent papers on code summarization
 - found two commonly used evaluation methodologies: mixed-project and cross-project
 - define two use cases that could be evaluated by these methodologies
- Define a more practical use case: continuous-mode
- Propose an appropriate evaluation methodology for this use case: time-segmented
- Experiment several existing ML models using the three methodologies

evaluation methodology

mixed-project (used by 15/18) cross-project (used by 4/18) time-segmented (proposed)

use case

in-project batch-mode cross-project batch-mode continuous-mode



Outline

• Evaluation methodologies and use cases

mixed-project	in-project batch-mode
cross-project	cross-project batch-mode
time-segmented	continuous-mode

• Experiments to study the impact of evaluation methodologies

- experiments setup
- dataset
- results and findings



Mixed-Project Evaluation Methodology & In-Project Batch-Mode Use Case

- Used in prior work
- Randomly shuffle the samples and split them into training, validation, and test sets





Cross-Project Evaluation Methodology & Cross-Project Batch-Mode Use Case

- Used in prior work
- Randomly shuffle the projects and split them into training, validation, and test sets



Training 🅕 Validation 🖨 Test





Limitation of Batch-Mode Use Cases

Usually happen only once in the lifecycle of a project





Continuous-Mode Use Case



Time-Segmented Evaluation Methodology

- Not used in prior work on developing new ML models for code summarization
- Split samples in a time-aware method
 - assign samples before τ^{-2} to training set
 - assign samples after τ^{-2} and before τ^{-1} to validation set
 - assign samples after τ^{-1} and before τ to test set





Outline

Evaluation methodologies and use cases

mixed-project	in-project batch-mode
cross-project	cross-project batch-mode
time-segmented	continuous-mode

• Experiments to study the impact of evaluation methodologies

- experiments setup
- dataset
- results and findings



Experiments Setup

Task	comment generation	method naming
Models	DeepComHybrid Hu et al. ESE'20 Transformer Seq2Seq	Code2VecAlon et al. POPL'19Code2SeqAlon et al. ICLR'19
Metrics	BLEU METEOR ROUGE-L EM (exact match)	Precision Recall F1 EM (exact match)



Dataset

- 77,745 (code, comment) with timestamps from 160 popular Java projects on GitHub
- Given a dataset of (code, comment) with timestamps, split it to get
 - training (Train), validation (Val), and standard test (TestS) sets for each methodology
 - common test (TestC) set to compare each pair of methodologies



Results and Findings (1/4)

- Different methodologies may lead to conflicting evaluation results
 - Code2Vec is better than Code2Seq under the mixed-project and time-segmented methodologies, but is worse under the cross-project methodology





Results and Findings (2/4)

- Different methodologies may lead to conflicting evaluation results
 - Transformer is statistically significantly better than Seq2Seq under the time-segmented methodology, but not under the cross-project methodology



*no significant difference between the models



Results and Findings (3/4)

- Evaluation results from prior work do not represent the ML models' performance in the continuous-mode use case
 - Results under the mixed-project methodology are inflated
 - Results under the cross-project methodology may be an under-estimation ٠



task: method naming



Results and Findings (4/4)

- Evaluation results from prior work do not represent the ML models' performance in the continuous-mode use case
 - Results under the mixed-project methodology are inflated
 - Results under the cross-project methodology may be an under-estimation







task: comment generation metric: BLEU

Conclusions

- We need to more diligently choose evaluation methodology and report results of ML models according to the intended use cases
- Time-segmented evaluation methodology should be adopted in the evaluation of ML models for code summarization

data and code: <u>https://github.com/EngineeringSoftware/time-segmented-evaluation</u> preprint: <u>https://arxiv.org/abs/2108.09619</u>

Pengyu Nie <pynie@utexas.edu>





backup slides

Evaluation Methodologies of ML Models

- Evaluation paradigm based on automatic metrics:
 - split a dataset of (code, comment) samples into training, validation, and test sets
 - train ML models on training + validation sets
 - report automatic metrics (e.g., BLEU, F1) on test set



- What is the intended use cases of the ML model?
- How to split the dataset, such that the evaluation results represent the ML model's performance in the intended use cases?
- In the context of code summarization: should we consider the timestamps of code and comments?

Timestamps of Code and Comments

- Developers iteratively add/edit code and comments
 - the style of newer code and comments written can be affected by older code and comments
- Temporal relations among samples are not explicitly modeled in the evaluation of prior work
 - can lead to inflated values for automatic metrics
 - can lead to misunderstanding if a model might be useful once adopted



Experiments Setup

- (code, comment) with timestamps from popular Java projects on GitHub using English for summaries
- collected samples before $\tau = 2021.1.1$ time-segmented on $\tau^{-2} = 2019.1.1$ and $\tau^{-1} = 2020.1.1$
- splitting ratios for in-project and cross-project: 70%,10%,20%
- Models
 - comment generation: DeepComHybrid, Transformer, Seq2Seq
 - method naming: Code2Vec, Code2Seq
- Automatic metrics
 - comment generation: BLEU, METEOR, ROUGE-L, EM (exact match)
 - method naming: precision, recall, F1, EM (exact match)
- Run each model 3 times & perform statistical significance tests

Task		Train	Val	TestS		TestC
ent	MP	50,879	7,569	14,956 5	$\longrightarrow \mathbf{MP} \cap \mathbf{CP}$	3,362
mme	CP	50,879	8,938	15,661 <	$ \longrightarrow \mathbf{MP} \cap \mathbf{T} $	2,013
Ger Co	T	50,879	11,312	9,870 4	\longrightarrow CP \cap T	2,220
ററ	MP	50,879	7,523	14,796 •	$\checkmark \mathbf{MP} \cap \mathbf{CP}$	3,344
etho amin	CP	50,879	8,811	15,332 <	$ \longrightarrow \mathbf{MP} \cap \mathbf{T} $	2,011
ΧŽ	T	50,879	11,223	9,807 4	\hookrightarrow CP \cap T	2,211

Results and Findings (1/3)



• Depending on the methodology, one model can perform better or worse than another

Results and Findings (2/3)

Train Test	$\frac{\mathbf{MP}}{\mathbf{MP}} \stackrel{\mathbf{CP}}{\cap \mathbf{CP}}$	P MP P MP	\mathbf{T} $\mathbf{T} \cap \mathbf{T}$	CP CP ∩	T T	
	V BL	EU [%]				
	45.0 11. 53.8 ^α 13. 58.1 ^α 14 .	4 ^β 52.6 7 61.7 1 65.6	43.2 ^β 53.2 56.3	11.0 13.4 14.3	45.6 53.3 56.1	
	▼ MET	EOR [%]				
DeepComHybrid Seq2Seq Transformer	52.9 18. 62.0 ^γ 21. 66.0 ^γ 21 .	3 659.3 2 68.0 4 71.5	50.0 ^δ 59.8 62.7	18.2 [¢] 21.2 [¢] 21.6	52.0 59.6 62.1	
		GE-L [%]				
	no 57.0 23. 66.7 29. 70.1 30.	t signifi 4 72.0 4 74.9	cant ⁵ 64.1 66.6		55.9 64.3 66.4	not significan
		M[%]				
	30.2 $\eta \theta$ 1.4 37.3 $\eta \iota$ 1.4 41 1 $\theta \iota$ 1	$\begin{array}{c c c} 4 & ^{\kappa}39.5 \\ 2 & 48.7 \\ 7 & 52.3 \\ \end{array}$	31.0 ^κ 41.0 44.2	$\lambda^{\mu} 1.3$ $\lambda^{1.1}$		

• Depending on the methodology, the differences between models may or may not be observable

Results and Findings (3/3)



- Results under the mixed-project methodology are inflated
- Results under the cross-project methodology may be an under-estimation of the more realistic continuous-mode use case

Conclusions

- We need to more diligently choose evaluation methodology and report results of ML models according to the intended use cases
- Time-segmented evaluation methodology should be adopted in the evaluation of ML models for code summarization
- Misuse of evaluation methodologies can lead to inflated values for automatic metrics and misunderstanding if a model might be useful once adopted

data and code: <u>https://github.com/EngineeringSoftware/time-segmented-evaluation</u> preprint: <u>https://arxiv.org/abs/2108.09619</u>

Pengyu Nie <pynie@utexas.edu>



Application of Methodologies: Goals

- Given a dataset of (code, comment) with timestamps, split it to get
 - training (Train), validation (Val), and standard test (TestS) sets for each methodology
 - common test (TestC) set to compare each pair of methodologies



MP = mixed-project CP = cross-project T = time-segmented

Application of Methodologies: Step 1/6

- Given a dataset of (code, comment) with timestamps, split it to get
 - training (Train), validation (Val), and standard test (TestS) sets for each methodology
 - common test (TestC) set to compare each pair of methodologies

1. time-segment samples in each project



Application of Methodologies: Step 2/6

- Given a dataset of (code, comment) with timestamps, split it to get
 - training (Train), validation (Val), and standard test (TestS) sets for each methodology
 - common test (TestC) set to compare each pair of methodologies



Application of Methodologies: Step 3/6

- Given a dataset of (code, comment) with timestamps, split it to get
 - training (Train), validation (Val), and standard test (TestS) sets for each methodology
 - common test (TestC) set to compare each pair of methodologies



Application of Methodologies: Step 4/6

- Given a dataset of (code, comment) with timestamps, split it to get
 - training (Train), validation (Val), and standard test (TestS) sets for each methodology
 - common test (TestC) set to compare each pair of methodologies



Application of Methodologies: Step 5/6

- Given a dataset of (code, comment) with timestamps, split it to get
 - training (Train), validation (Val), and standard test (TestS) sets for each methodology
 - common test (TestC) set to compare each pair of methodologies



MP = mixed-project CP = cross-project T = time-segmented

Application of Methodologies: Step 6/6

- Given a dataset of (code, comment) with timestamps, split it to get
 - training (Train), validation (Val), and standard test (TestS) sets for each methodology
 - common test (TestC) set to compare each pair of methodologies



Mixed-Project Evaluation Methodology

Randomly shuffle the samples and split them into training, validation, and test sets



In-Project Batch-Mode Use Case



Cross-Project Evaluation Methodology

Randomly shuffle the projects and split them into training, validation, and test sets



Cross-Project Batch-Mode Use Case





do not write comments for any method

time to add documentations, with ML model



time↓

τ

Not Considering Temporal Relations During Evaluation

- Temporal relations among samples are not explicitly modeled
- Can lead to misunderstanding if a model might be useful / not useful once adopted



time

Experiments Setup

Task comment generation	method naming
-------------------------	---------------