

# Going beyond “Good Job”: Analyzing Helpful Feedback from the Student’s Perspective.

M Parvez Rashid, Yunkai Xiao, Edward F. Gehringer  
Department of Computer Science  
North Carolina State University  
Raleigh, NC, USA  
{mrashid4, yxiao28, efg}@ncsu.edu

## ABSTRACT

Peer assessment can be a more effective pedagogical method when reviewers provide quality feedback. But what makes feedback helpful to reviewees? Other studies have identified quality feedback as focusing on detecting problems, providing suggestions, or pointing out where changes need to be made. However, it is important to seek students’ perspectives on what makes a review helpful to a reviewee. This study explores the helpfulness of feedback from students’ perspectives when the feedback contained suggestions or mentioned problems or both. We applied natural language processing techniques to identify suggestions and problems mentioned in peer reviews. We also analyzed important text features that are associated with suggestions or problems detected by the peer feedback. The result showed that students are likely to find a review helpful if a suggestion is provided along with the problem mentioned in the feedback rather than simply identifying the problem.

## Keywords

Peer assessment, neural-network, natural language processing, correlation coefficient, suggestions

## 1. INTRODUCTION

Peer assessment has been proven to be an effective learning approach in both face-to-face and distance learning classes. It is especially useful in massive open online courses (MOOCs) where the potentially overwhelming number of students has no fixed bound. All of these students must be assessed by someone, and there are only a limited number of staff. Peer assessment can be as accurate instructor assessment, since artifacts are reviewed by multiple assessors who can invest more time than a teacher could [17]. It can also provide timely feedback [1] that helps students to focus on their weaknesses. Peer assessors pick up some of the feedback workload for instructors, who can then offer more help to students who are in need.

M. P. Rashid, Y. Xiao, and E. F. Gehringer. Going beyond “good job”: Analyzing helpful feedback from the student’s perspective. In A. Mitrovic and N. Bosch, editors, *Proceedings of the 15th International Conference on Educational Data Mining*, pages 515–521, Durham, United Kingdom, July 2022. International Educational Data Mining Society.

© 2022 Copyright is held by the author(s). This work is distributed under the Creative Commons Attribution NonCommercial NoDerivatives 4.0 International (CC BY-NC-ND 4.0) license.  
<https://doi.org/10.5281/zenodo.6853179>

Peer assessors provide assessment in two forms. One is textual feedback, which takes the form of prose feedback to a peer. This is usually used as formative feedback. Another is numerical scores, on a Likert scale, which allows a summative grade to be calculated. Most online peer-assessment environments use both kinds of feedback. Studies have shown that students learn more from giving feedback than receiving it [13, 2, 7, 5] and giving feedback engage students in active learning [16]. It forces students to think metacognitively [4], and learn in-depth, as reviewing a peer requires a good hold on the topic [8].

However, the learning experience in a peer-review environment depends on the quality of reviews provided by the student peers. The goals of fairness and equity require that, insofar as possible, all students receive helpful formative feedback on their work. But, not all assessors provide constructive feedback, due to lack of knowledge in the topic or simply carelessness. To encourage and guide students in reviewing the artifacts, instructors typically need to scrutinize reviews manually. This consumes a good portion of the time that would be saved by having students provide quality feedback. An automated analysis could save considerable time.

A few studies [20, 19] have tried to lessen the instructors’ assessment burden by automatically detecting characteristics of a quality review. That raises the question of what defines a quality review. According to Nelson and Schunn [10] high-quality feedback consists of (i) identifying a problem and (ii) suggesting a solution. However, their finding was based on students’ performance and not from their (students’) perspective. It is important to identify whether “quality feedback” is actually helpful to the reviewees, based on students’ opinion of what feedback is helpful.

In this paper, we propose a method using natural language processing (NLP) and neural networks to automate the process of analyzing and classifying reviews to discover whether they contain suggestions and/or problems. We analyzed the words that are used to include suggestions or problems in feedback. Our goal is to answer the following research questions:

- **RQ1:** Can we build a model to accurately detect comments containing suggestions or detecting problems?
- **RQ2:** Are “quality comments”—those containing suggestions, detecting problems, or both—actually helpful from the student’s point of view?

- **RQ3: Can an automated process effectively identify helpful feedback?**

## 2. RELATED WORK

This section discusses related work on identifying the properties of a “quality“ or “helpful“ peer assessment.

Nelson and Schunn [10] examined five features of feedback (summarization, specificity, explanations, scope, and affective language) that constitute good-quality reviews, and the correlations among them. Their study divided the features of feedback into cognitive and affective components. According to their findings, summarization, specificity, explanations, and scope are cognitive in nature. Cognitive features of a review are expected to most strongly affect understanding. This explanation helped us to identify suggestions and problem detection as a property of quality feedback.

An approach to improve review quality is to provide the reviewer with a rubric defining the characteristics of a quality review. Jaco du Toit [3] conducted a study to identify the impact of peer review on essay assignments. The study showed that giving students a rubric describing the characteristics of a good essay can provide them with the insight to produce better quality assessments than they would otherwise produce. But this study did not specify the qualities of a good review. When they received poor reviews, they were confused about the quality of their work, sometimes feeling a false sense of accomplishment. Rashid et al [15] analyzed rubric items to determine which of them induce peer reviewers to write quality feedback using NLP approaches. In their work quality feedback was identified if the review text contained a suggestion, detected a problem problem, or was localized (pinpointing the place where a revision should be made).

McGrath and Taylor did a study on students’ perception of helpful feedback for writing performance [9]. Their study defined quality feedback (“developed feedback”) as clear, specific, and explanatory in nature. They measured students’ perception of developed feedback by having them rate the feedback on a Likert scale. The results showed that students rated developed feedback highly for helpfulness.

A survey of 44 students done by Weaver showed that, in order to use the feedback, students needed advice (suggestions)[18]. The analysis of the feedback content and students’ responses uncovered that vague feedback (e.g., “Good job”) is unhelpful, lacking in guidance (void of suggestions), or focused on the negative (mentioning only problems), or was unrelated to assessment criteria.

Ramachandran et al. [14], developed an automated system to evaluate reviews and show how they compared to other reviews for the same assignment. They extracted attributes like relevance to the submission, content, coverage, tone, and volume of feedback to identify a good-quality review. They constructed word-order graphs to compare the reviews with submission text and extract features from the reviews.

To identify localization and make suggestions to improve the review Nguyen et al. [11] applied natural language processing techniques. They provided real-time formative feedback

to reviewers on how to localize their review comments.

Zingle et al. [20] used neural-network approaches to find suggestions in the review text, and compared them against rule-based NLP approaches. In a similar work Xiao et al. [19] used NLP techniques with several ML and neural-network approaches to identify problem statements in review text. Our work takes this a step further and asks whether it is enough for a review to detect problems, or whether reviews that also make suggestions are more helpful.

## 3. DATA

Machine-learning and neural network-based models can perform as well or as badly as the data they are given. However, obtaining good labeled data is expensive. For the purpose of our experiment, we have collected labeled datasets for comments with three different characteristics:

- detects a problem: A review comment is labeled yes or no according to whether it detects a problem.
- contains a suggestion: A review comment is labeled yes or no according to whether it contains a suggestion.
- is helpful: A review comment is labeled yes or no depending on whether the reviewee found it helpful.

We acquired this labeled peer-review data from the Expertiza system in a systematic manner. Expertiza is a system to support different kinds of communications that are involved in the peer-assessment process. It supports double-blind communications between authors and reviewers, assessment of teammate contributions, and evaluations by course staff.

For the purpose of this study, we collected the data from Object-Oriented Design and Development course at NC State University for about three years. This course used the Expertiza system to manage the peer-review assessment process for evaluating the students. In each semester, this course typically assigns three peer-reviewed assignments to students, who work in teams consisting of two to four members. Even though the assignments are done in a group setting, the submissions are reviewed by individual students from other groups. After receiving the reviews from peers, teams revise their work and resubmit it for grading. The second round of the assessment is generally summative, where along with textual comments, the peer-reviewers assign scores to the submission.

Generally, a small number of people cannot annotate a large dataset. It is better to have a large number of people each undertake a small number of annotation tasks; this lessens the chance that an annotator will become fatigued and assign inaccurate labels. We engaged students in the labeling task by offering a small amount of extra credit. After receiving peer feedback, students were asked to label the feedback to identify whether the reviewer mentioned a problem or suggested a solution. They were also asked whether they considered the feedback to be helpful. In different assignments, students were asked to label the feedback for different characteristics; the same comments were not necessarily labeled for all three characteristics. After labeling was complete, the course instructor and TAs spot-checked the data

**Table 1: Sample review comment and annotations done by students ('1' indicates 'yes' and '0' indicates 'no')**

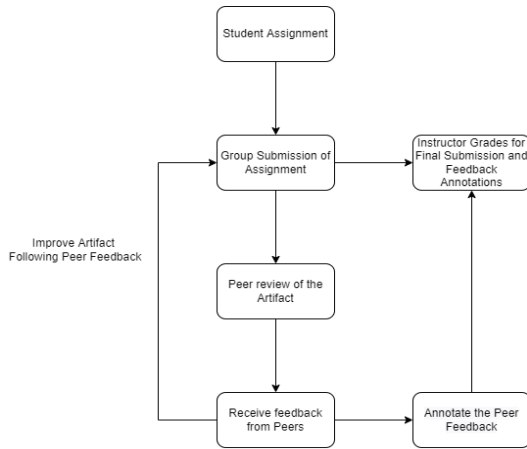
Review Comment	Detects Problem
The Travis CI Build is Failing as of now. No conflicts as per the GitHub report.	1
Yes, the explanation is elaborative and complete.	0
Since the build failed, I would not recommend adding it to the production server yet.	1

Review Comment	Gives Suggestion
Test Plan is too verbose. Trivial areas can be trimmed off.	1
The team needs to look into Travis CI log & 1	1
Many test cases in terms of controllers, but none for models.	0

Review Comment	Is Helpful
The build is failing due to 4 failures in the model specs.	1
The writeup is clear.	0
Since the build failed, I would not recommend adding it to the production server yet.	1



**Figure 1: Flow diagram of peer review and feedback annotation process**

that each student labeled. If any labels were found to be incorrect, the data labeled by that student was excluded from the dataset.

Since the reviews were done on team projects, and labeling was done individually, two to four students had the opportunity to label (or “tag”) the same review comments. If multiple students did tag the same comment, inter-rater reliability (IRR) could be calculated. We chose Krippendorff’s  $\alpha$  [6] as the metric for IRR. We chose this metric because it is not impacted by missing ratings, which were common since not all students availed themselves of the extra-credit opportunity. In an effort to use only the most reliable labeling, we included only labels that were assigned (or not assigned) by all the students in the team that was reviewed. This allowed us to raise Krippendorff’s  $\alpha$  of our dataset from 0.696 to 1. Figure 1 shows the peer-review and annotation process.

Following the described process we accumulated 18,392 annotations for problem detection, 7,416 for suggestion-detection and 3,970 for helpfulness-detection datasets. All the three datasets have an equal ratio of the binary class labels (i.e., they are balanced). Sample comments from the three datasets are shown in Table 1

## 4. METHOD

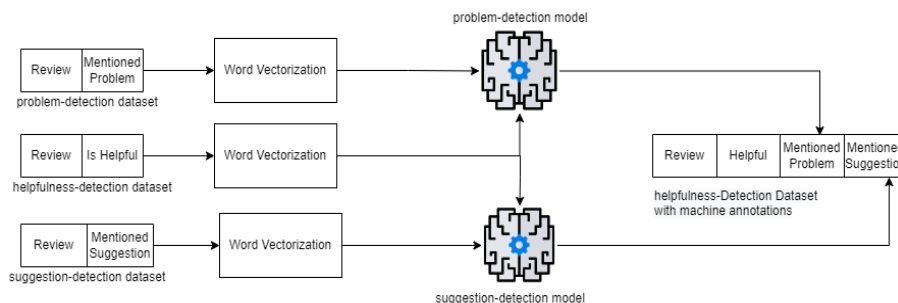
Our goal in this study is to analyze students’ perspectives on helpful comments that mentioned problems and/or suggestions. To conduct the study, we had students annotate comments on the basis of whether they found them helpful. We need an automated process to identify those review comments that contain suggestions and/or problem statements. We first train a model (the problem-detection model) to classify reviews that contain a problem statement by training and testing with the problem-detection dataset. We build a second model (the suggestion-detection model) to classify the presence of suggestions in a review comment by using the suggestion-detection dataset. As model performance matters, we applied several ML and neural-network models to pick the most accurate models for annotating the helpfulness-detection dataset. Figure 2 shows the annotation process of the helpfulness-dataset using the models.

When approaching a classification problem by any type of machine-learning (ML) or neural network models, there are many different approaches to choose from. No one model is best for all problems. In our study, we have chosen Support Vector Machine (SVM), Random forest (RF), classical ML models and compared their performance with Bi Directional Long Short-term Memory (Bi-LSTM), and Bidirectional Encoder Representations from Transformers (BERT) models. We used TF-IDF for ML models and Global Vectors for Word Representation (GloVe) for Bi-LSTM to perform word vectorization. Before we applied any word vectorization techniques, we cleaned the text by removing URLs, stop words, and applying stemming. We use our problem-detection dataset and suggestion-detection dataset on these model with 80:10:10 ratio for training, testing and validation.

### 4.1 Classical Machine-Learning Models

#### 4.1.1 Input Embedding with TF-IDF

Machine-learning models are suitable for capturing complex relationships between the input data. But they require numeric input. The review data that we have in our dataset is textual. We have to convert them to numbers and also allow the model to capture the important features of the text. One way to do that is term frequency-inverse document frequency (TF-IDF). TF-IDF measures the importance of a word in a document using statistical calculation. If a word appears more times in a document the importance of the word in the document increases proportionally. We used



**Figure 2: Annotation process of the helpfulness-dataset for mentioned problem and suggestions in the comments using models. The models were trained for detecting problems or suggestions mentioned in the review text. The training datasets were annotated by human (students).**

scikit-learn [12] library to implement TF-IDF and vectorize the words in the feedback.

### 4.1.2 SVM

SVM is very popular for high accuracy and low computational cost. For a classification problem between two classes, SVM maximizes the margin of the separation plane between the two classes. We provided the feature vector of the reviews converted by TF-IDF to the SVM model to classify the review for having a particular property (contains problem or suggestion in the comment). We applied a grid search to find the best inverse regularization parameter  $C$ .

### 4.1.3 RF

We used Random Forest for its popularity to make more accurate classification with a simple approach. RF makes an ensemble decision from a forest consisting of multiple uncorrelated decision trees. The general idea of the RF is that the decision from individual decision trees increase the accuracy of overall result. We varied the number of decision trees and depth of the trees to get the best result. We used TF-IDF for making feature vectors from the review text.

## 4.2 Neural-Network Models

### 4.2.1 Input Embedding

Neural network models are popular for text classification tasks. However, to improve the performance of the neural-network models on the text data, it is necessary to represent the data that is suitable for the model to work with, and without losing the underlying latent relations among the features of the data. For our experiment we have used GloVe with Bi-LSTM. GloVe not only measures the statistical significance of words, it also considers the statistical co-occurrence and semantic relation of the words.

### 4.2.2 Bi-LSTM

Bi-Long Short-Term Memory is in general used for sequential data classification tasks. It is a good fit for peer-review texts. Review comments are sequential data, and the words of the text have latent semantic and contextual relations with each other. As Bi-LSTM model takes input from both right and left direction of the text, it can capture the relationship between the words in texts occurring in any order.

**Table 2: Hyperparameters of Models**

Model	Hyperparameter
SVM	$c=1$
RF	tree = 100 max depth = 4
Bi-LSTM	maximum text length = 300 Embedding = 300d Hidden layer activation = ReLu dropout = 0.4 optimizer = Adam Output layer activation = Sigmoid Epoch=20
BERT	optimizer = AdamW Learning rate = $2e-5$ Epoch=4

### 4.2.3 BERT

BERT is based on Transformer model and use attention mechanism to learn the contextual relations of the words in a sentence. Being a bi-directional input reader, BERT learns the context of word in sentence by considering words occurring before and after.

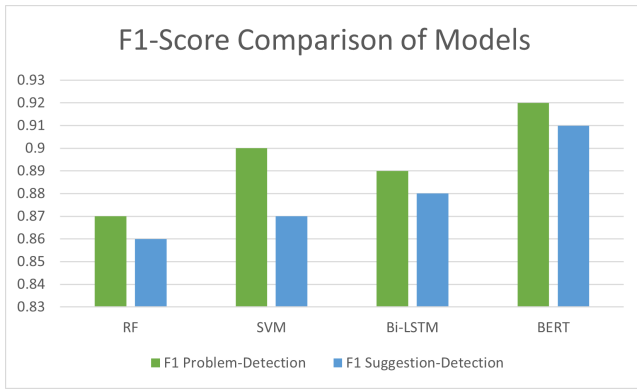
## 5. RESULTS

In this study, if a feedback comment mentions problems and/or suggestions, we are considering it to be quality feedback. Our first step is to construct two separate models where one identifies whether feedback contains a problem statement and another identifies whether feedback contains a suggestion. To identify the best-performing models we trained and tested the performance of several classical ML models and neural-network models and compared their performance.

**RQ1: Can we build a model to accurately detect comments containing suggestions or detecting problems?**

Figure 3 reports the comparison of the F1-score values of the classical machine-learning (ML) models and neural-network models on the problem-detection dataset and suggestion-detection dataset. To compare the performance of the models we use the F1-score, as this represents the harmonic mean of precision and recall.

- **On the problem-detection dataset:** Among the classical



**Figure 3: F1-score comparison to measure performance on classifying review text on problem detection and suggestion detection using classical ML and neural-network models. In overall F1-score comparison, the BERT model shows the best performance.**

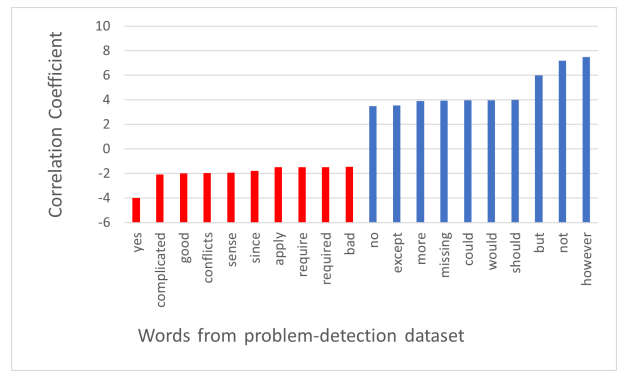
ML models SVM made the highest f1-score 0.90 and among the neural-network models; BERT obtains the overall highest F1-score, 0.92.

- **On the suggestion-detection dataset:** BERT achieved the highest F1-score, 0.91. Among the classical ML models, SVM achieved the highest F1-score, 0.87.

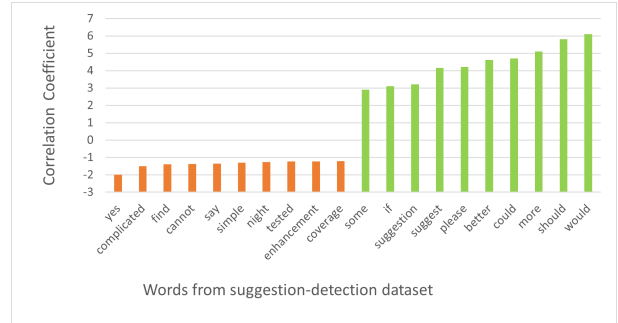
To gain a deeper insight into the words that are highly correlated with text where a problem mention or suggestion is mentioned, we analyzed the top 10 positive and negative correlation coefficient values calculated by the SVM model. Figure 4(a) shows the coefficient values that the problem-detection model has calculated for various words. Note that it has positive coefficient values for words such as “however”, “but”, and “not”. In the English language these words are more likely to be used when stating problems. Similarly words like “yes”, “completed” and “good” are not likely to occur in a problem statement. Figure 4(b) shows that the suggestion-detection model has positive coefficient values for words like “should”, “would”, “more”, “suggest”. These words are likely to be used in suggestions. On the other hand, words “yes”, “completed”, and “cannot” are more likely not to be used to express suggestions; thus they have negative coefficient values.

As BERT outperformed all other models on both problem and suggestion datasets, we trained two separate BERT models to annotate the feedback comments contained in the helpfulness-detection dataset. The BERT-created annotations recorded whether each comment in the helpfulness-detection dataset detected a problem or offered a suggestion. The models annotated each comment with either “1” or “0”, indicating having the property or not. We perform an **and**-operation using the BERT-created annotations. If both the problem and suggestion were mentioned in a comment the **and**-operation yields 1 otherwise 0. The resulting helpfulness-detection dataset is shown in Table 3.

**RQ2: Are “quality comments”—those containing suggestions, detecting problems, or both—actually helpful from the stu-**

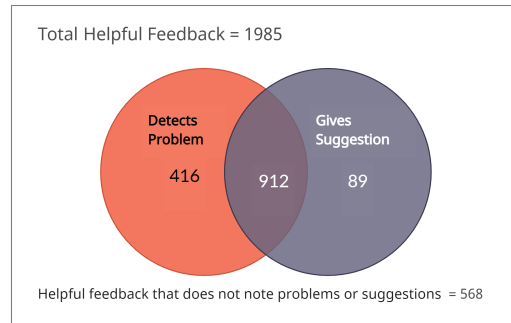


(a)

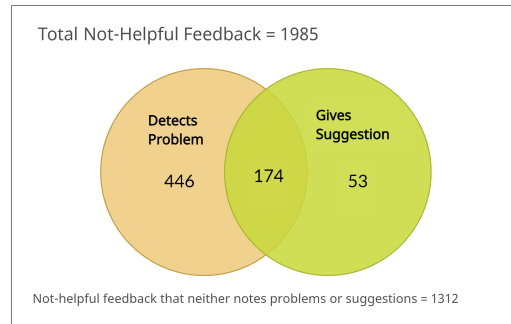


(b)

**Figure 4: Top 10 positive and negative coefficient value of words from the problem-detection and suggestion-detection datasets**



(c)



(d)

**Figure 5: Venn diagram of helpful feedback annotated for mentioned suggestion and/or problem**

**Table 3: Table shows sample comments from helpfulness-detection dataset and corresponding annotations.** Note that “Is Helpful” annotations are done by humans (students), while “Detects Problem” and “Gives Suggestion” are annotated by the BERT model. “Contains Problem and Suggestion” is from **anding** the “Detects Problem” and “Gives Suggestion” columns.

Review Comment	Is Helpful (human-annotated)	Detects Problem (machine-annotated)	Gives Suggestion (machine-annotated)	Contains Problem and Suggestion (and-operation)
The build is failing due to 4 failures in the model specs.	1	1	0	0
The writeup is clear.	0	0	0	0
Since the build failed, I would not recommend adding it to the production server yet.	1	1	1	1
I would recommend adding more code for helping following their changes.	1	0	1	0

### dent’s point of view?

After we computed the annotations for problem detection and suggestions, we did a Venn diagram analysis on the updated helpfulness-detection dataset. The diagrams illustrate the overlap of comments that both detect a problem and offer a suggestion. Figure 5(c) shows that 1,985 comments in the helpfulness-detection dataset were annotated by students as being helpful. Among the helpful comments, 1,417 were annotated for having problems and/or suggestions mentioned in the feedback. Out of these 1,417 helpful comments, 912 of them were machine-annotated as containing both problem detection and suggestions. A total of 568 helpful comments did not have any problem or suggestion mentioned, based on machine-annotation.

On the other hand, out of the 1,985 comments that were human-annotated as not helpful [Figure 5(d)], 673 comments were annotated as having either a problem and/or suggestion mentioned. Among those 673 comments, only 174 were annotated as having both suggestion and problem mentioned. A total of 1,312 comments that did not have any problem or suggestion mentioned were annotated as not helpful by the students.

To summarize the Venn diagram analysis, comments that the students found helpful mostly detected problems and/or contained suggestions. However, among those comments noting suggestions and/or problems, students annotated as helpful mostly comments that *both* pointed out problems and gave suggestions. This indicates that peer feedback is more helpful to the students when a suggestion is given in a comment that detects a problem. On the flip side, Figure 5(b) suggests that students rarely find comments helpful when they do not mention any problem or contain a suggestion.

### RQ3: Can an automated process effectively identify helpful feedback?

A key question is whether comments automatically annotated as “quality” (meaning that they both identified a problem and gave a suggestion) were the same comments that the students considered helpful (that they manually labeled as helpful). Among the comments that the students considered helpful, 64% of them both mentioned a problem and gave a suggestion. Conversely, of the comments that the students labeled as not helpful, 66% of them neither mentioned a problem nor contained a suggestion.

The results indicate that the automated annotation per-

formed by the BERT model can be very effective in predicting which comments students will consider helpful. While it can’t deliver an actual count of helpful comments in a particular review, that is not important. It *can* determine whether the feedback provided by the reviewer contains a substantial number of quality comments. That is what is needed to automatically detect helpful reviews.

## 6. CONCLUSION

This study constitutes the first analysis of the helpfulness of peer-assessment feedback from student’s perspective. Feedback that mentions problems or includes suggested changes was considered to be quality feedback. We used natural language processing (NLP) techniques in conjunction with several ML and neural networks to identify quality peer feedback. We systematically collected and scrutinized 18,392 comments mentioning problems, 7,416 comments containing suggestions, and 3,970 comments that were annotated by humans (students) as being helpful.

Using the annotated dataset, we trained our ML and neural-network models to identify quality feedback. For identifying suggestions and problems mentioned in the review text, the BERT model outperformed the other models. As the BERT model focuses on the important features of the text, it was best at identifying suggestions and problems in the feedback. We used the BERT model to automatically annotate comments as mentioning problems or making suggestions, and compared these annotations with comments that students had manually annotated as being helpful. We also analyzed important words that are frequently present in comments mentioned a problem or suggestion.

A key finding of this study is that the students find review comments more helpful when peer reviewers both mention problems in the reviewed artifact and provide suggestions on how to resolve the issue. We can use a state-of-the-art BERT model to automatically identify the helpful review comments.

It should not be hit-or-miss whether students receive helpful reviews on their submitted work from peer reviewers. This study helps identify helpful feedback and therefore, help students to improve their work.

## 7. REFERENCES

- [1] J. Cambre, S. Klemmer, and C. Kulkarni. Juxtapeer: Comparative peer review yields higher quality feedback and promotes deeper reflection. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2018.
- [2] Y. D. Çevik. Assessor or assessee? investigating the differential effects of online peer assessment roles in the development of students’ problem-solving skills. *Computers in Human Behavior*, 52:250–258, 2015.
- [3] J. du Toit. Enhancing the quality of essays through a student peer-review process. In *International Conference on Innovative Technologies and Learning*, pages 459–468. Springer, 2019.
- [4] E. F. Gehringer. A survey of methods for improving review quality. In *International Conference on Web-Based Learning*, pages 92–97. Springer, 2014.
- [5] M. H. Graner. Revision workshops: An alternative to peer editing groups. *The English Journal*, 76(3):40–45, 1987.
- [6] K. Krippendorff. *Content analysis: An introduction to its methodology*. Sage publications, 2018.
- [7] L. Li and V. Grion. The power of giving feedback and receiving feedback in peer assessment. *All Ireland Journal of Higher Education*, 11(2), 2019.
- [8] K. Lundstrom and W. Baker. To give is better than to receive: The benefits of peer review to the reviewer’s own writing. *Journal of second language writing*, 18(1):30–43, 2009.
- [9] A. L. McGrath, A. Taylor, and T. A. Pychyl. Writing helpful feedback: The influence of feedback type on students’ perceptions and writing performance. *Canadian Journal for the Scholarship of Teaching and Learning*, 2(2):5, 2011.
- [10] M. M. Nelson and C. D. Schunn. The nature of feedback: How different types of peer feedback affect writing performance. *Instructional Science*, 37(4):375–401, 2009.
- [11] H. Nguyen, W. Xiong, and D. Litman. Iterative design and classroom evaluation of automated formative feedback for improving peer feedback localization. *International Journal of Artificial Intelligence in Education*, 27(3):582–622, 2017.
- [12] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- [13] R. Rada et al. Collaborative hypermedia in a classroom setting. *Journal of Educational Multimedia and Hypermedia*, 3(1):21–36, 1994.
- [14] L. Ramachandran, E. F. Gehringer, and R. K. Yadav. Automated assessment of the quality of peer reviews using natural language processing techniques. *International Journal of Artificial Intelligence in Education*, 27(3):534–581, 2017.
- [15] M. P. Rashid, E. F. Gehringer, M. Young, D. Doshi, Q. Jia, and Y. Xiao. Peer assessment rubric analyzer: An nlp approach to analyzing rubric items for better peer-review. In *2021 19th International Conference on Information Technology Based Higher Education and Training (ITHET)*, pages 1–9. IEEE, 2021.
- [16] K. Topping. Peer assessment between students in colleges and universities. *Review of educational Research*, 68(3):249–276, 1998.
- [17] K. J. Topping. Peer assessment. *Theory into practice*, 48(1):20–27, 2009.
- [18] M. R. Weaver. Do students value feedback? student perceptions of tutors’ written responses. *Assessment & Evaluation in Higher Education*, 31(3):379–394, 2006.
- [19] Y. Xiao, G. Zingle, Q. Jia, H. R. Shah, Y. Zhang, T. Li, M. Karovaliya, W. Zhao, Y. Song, J. Ji, et al. Detecting problem statements in peer assessments. In *Proceedings of The 13th International Conference on Educational Data Mining (EDM 2020)*, pages 704–709, 2020.
- [20] G. Zingle, B. Radhakrishnan, Y. Xiao, E. Gehringer, Z. Xiao, F. Pramudianto, G. Khurana, and A. Arnav. Detecting suggestions in peer assessments. *International Educational Data Mining Society*, 2019.