

# How Teachers Affect Students' Online Participation in EFL Courses in Uruguay

Cecilia Aguerrebere<sup>1</sup>      Monica Bulger<sup>2</sup>      Cristóbal Cobo<sup>1</sup>

Sofía García<sup>1</sup>      Gabriela Kaplan<sup>3</sup>      Jacob Whitehill<sup>4</sup>

<sup>1</sup>Fundación Ceibal, Uruguay

<sup>2</sup>Future of Privacy Forum

<sup>3</sup>Plan Ceibal, Uruguay

<sup>4</sup>Worcester Polytechnic Institute, USA

April 28, 2020

## **Abstract**

We analyze teachers' written feedback to students in an online learning environment, specifically a setting in which high school students in Uruguay are learning English as a foreign language with both a classroom teacher and a remote teacher. We explored which factors are associated with greater student participation. How complex should teachers' feedback be? Should it be adapted to each student's English proficiency level? How does teacher feedback affect the probability of engaging the student in a conversation? We conducted both parametric multilevel modeling and non-parametric bootstrapping analyses of 27,627 messages exchanged between 35 teachers and 1074 students in 2017 and 2018. Our results suggest: (1) Teachers should

adapt their feedback complexity to their students' English proficiency level. Students who receive feedback that is too complex or too basic for their level post 13-15% fewer comments than those who receive adapted feedback. (2) Feedback that includes a question is associated with higher odds-ratio (17.5-19) of engaging the student in conversation. (3) For students with low English proficiency, slow turnaround (feedback after 1 week) reduces this odds ratio by 0.7. (4) For distance learning contexts such as this one, the classroom English teachers (CTs) – who both teach students locally and promote students' participation in the online program – may significantly affect students' commenting behavior. These results have potential implications for online platforms offering foreign language learning services, in which it is crucial to give the best possible learning experience while judiciously allocating teachers' time.

## 1 Introduction

For decades, teacher feedback has been shown to be one of the greatest drivers of student learning (Hattie and Timperley, 2007). The research focus has shifted from assessing whether feedback is effective to identifying the most effective strategies (Van der Kleij et al., 2015). Because of the complexity of the feedback process, the answer to this question remains deeply tied to the particular context in which it takes place. One particular learning domain that is growing fast in terms of number of learners and learning platforms (e.g., Duolingo, Babbel, Learning English at Coursera) is online learning of English as a foreign language (EFL). Despite its increasing prominence, teacher feedback in the online EFL context has received limited research attention (Van der Kleij et al., 2015; Conrad and Dabbagh, 2015).

In this paper we seek to contribute to the understanding of how teacher feedback influences students' behavior in the online EFL context. In particular, we focus on an EFL program in which students learn English with the help of both a local classroom English teacher (CT), who promotes students' participation in the program but does not provide

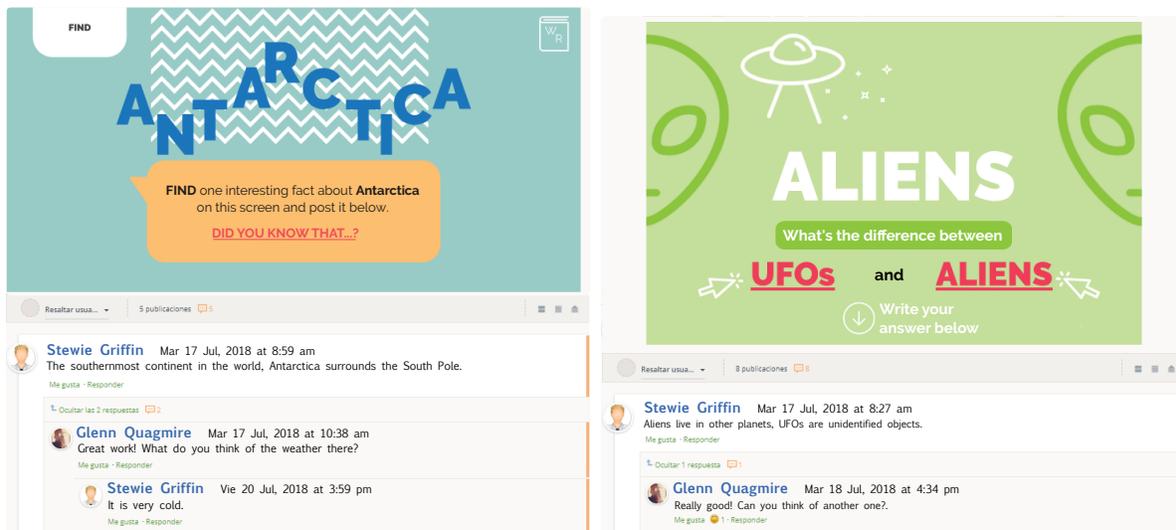


Figure 1: Example of a TDL exercise available through the LMS: after reading a material about Antarctica (left) or Aliens (right), the students are asked to mention interesting facts about the topic.

online feedback, and a remote teacher (RT), who is a native English speaker with whom students communicate online using discussion forums. Within this context, we seek to build an understanding of how the feedback the RTs give to their students affects their posting behavior: (1) How complex should the RT feedback be? (2) Should it be somehow adapted to their student's English proficiency level? (3) How does RT feedback affect the probability of engaging the student in a conversation? (4) What impact on students' behavior does the CT have? This research has potential implications for the countless online platforms offering foreign language learning services, aiming to enhance students' learning experience. This paper also makes a modest methodological contribution by illustrating how to conduct bootstrap analyses for data collected with nested structure; we hope these may inspire similar analyses for related contexts.

**Learning context:** This study is conducted in the context of a program for EFL learning created for secondary school students who attend the public school system in Uruguay. Uruguayan secondary school students (native Spanish speakers) often struggle with English,

having very disparate proficiency levels when they enter high school. This program, known as Tutorials for Differentiated Learning (TDL), was conceived to help tackle this problem by providing students with the option to learn and practice English at their own pace. For this purpose, a set of resources and exercises for EFL learning are made available online through an LMS system (see Figure 1). Completing an exercise consists of reading the material and posting a comment in English in a discussion forum. Exercises are organized into topics (e.g., music, sports, fashion, national parks, travel, etc), and there is one discussion forum per exercise. An RT, assigned to each classroom, reviews the students' posts and gives them individualized feedback. Note that, in contrast to RTs, the CT interacts with the students locally in the classroom twice a week during the English course. The students may be encouraged by their CT to explore the material and complete the exercises, but participation in the program is not mandatory. In the TDL learning context, CTs are encouraged to serve the role of program *promoters*, who foster student participation in TDL and help keep them on-task and engaged. This is somewhat different from other distance learning programs from the past decades (Alshammari et al., 2017; Kimball, 2002; Berge and Collins, 1995; Hannum et al., 2008); in those programs, CTs do not teach the subject matter themselves but rather serve as a facilitator between the pupils and the RT.

**RT-student interactions:** The student always starts the thread by posting a comment about a given topic in the LMS discussion forum. The RT replies giving the students personalized feedback on what they wrote. Then, the conversation may or may not continue depending on whether the student posts a new comment in the given thread. If the student does not, then that conversation ends there, but the student may start new threads when doing new exercises. Here is an example of an interaction where the student continued the conversation with the RT:

**Student:** *I do not have favorite music I like to listen to everything a little.*

**RT:** *That's great Alicia. What's your favourite song right now?*

**Student:** *At this moment I've heard a song from Michael Jackson that I loved its name is Thriller.*

**RT:** *Ok Alicia, thank you for sharing that :)*

and another example where she did not:

**Student:** *I see six oceans: Atlantic, Indian, Pacific, Atlantic, Arctic and Southern ocean.*

**RT:** *Very well Andrea.*

The TDL targets students with very diverse English proficiency levels, mostly quite low. The RTs' feedback is intended less as a way of correcting students' mistakes and more as a way to encourage students to participate and helping them realize that they are capable of communicating in English. Participation in the discussion forums is expected to conduce better learning since doing the exercises requires *reading* the material in English as well as *writing* the response in English. Therefore, two measures of interest are: the total comments the student posts, and whether the student engages in a given conversation with the RT.

## 2 Previous Work

Even though there is no unified definition of feedback, the seminal work by Hattie and Timperley (2007) conceptualizes feedback as *information provided by an agent regarding aspects of a student's performance or understanding*. It can be provided effectively, but it is dependent on several factors such as the task, the learning context, and the learners (Hattie and Timperley, 2007; Van der Kleij et al., 2015). It may improve learning outcomes when it has a direct use (e.g. correct the task), or it may increase motivation when only expressing praise for the student (Van der Kleij et al., 2015).

In the online language learning context, feedback has been reported as a fundamental aspect in skills development (Kahraman and Yalvac, 2015). Teacher feedback in online language learning environments can also inform development of data-driven personalized feedback. Emerging data-driven learning systems adapt feedback to individual student needs,

and have been shown to improve learning outcomes (Romero and Ventura, 2013). Furthermore, data mining has been used to understand how polarity (positive vs. negative comments) and timing can affect students' learning (Lang et al., 2015; Olsen et al., 2015).

Online teacher feedback, delivered asynchronously as in this study, offers students flexibility to read it at a time convenient to them, gives them the opportunity to concentrate more thoroughly on the comments in the absence of their peers, and to consult it whenever they complete future assessments (Hepplestone et al., 2011). Shang (Shang, 2017) presents a comparison of asynchronous peer feedback and synchronous corrective feedback in an online EFL environment for 44 university students in Taiwan. The major findings suggest that, even if participants accepted both approaches and obtained satisfactory results, they preferred the asynchronous modality as they consider peer discussions to be essential for clarifying misunderstandings or aspects of their writing.

Research on feedback for EFL learning in computer-mediated (CM) environments has widely focused on *peer* feedback, often on EFL writing (Guardado and Shi, 2007; Ho, 2015; Saeed et al., 2018). Jiang and Ribeiro (2017) present a systematic literature review on the effect of CM peer-written feedback on adult EFL writing. They confirmed the findings from previous research acknowledging the positive impact of CM peer feedback in this context. Nevertheless, the authors identified a series of factors on which the results are conditioned, such as the CM technology used, the types and content of peer feedback, as well as the learners' English proficiency level and technology-use anxiety.

As in many other subjects in educational data mining, most research on feedback has focused on higher education settings (Van der Kleij et al., 2015; Dekhinet, 2008), leaving the primary and secondary education contexts largely unexplored. We find previous work on secondary and primary education contexts on particular topics such as teachers' feedback strategies (MacDonald, 2015; Baadte and Kurenbach, 2017), student-generated feedback (Harris et al., 2015) among other more general examples (Oinas et al., 2017).

Related to our analyses of language complexity is a set of studies by Sinclair and colleagues. They studied the alignment in language complexity between a student and his/her conversation partner, which could be either another human (Sinclair et al. (2017)) or an automated conversation agent (Sinclair et al. (2019)). They found evidence that students and teachers adapt the complexity of their language based on the complexity they encounter in their conversation partners (Sinclair et al. (2017)).

Our work complements and enriches the previous work in several aspects: (1) it studies asynchronous teacher feedback in an online EFL environment, which has been seldom studied (Hepplestone et al., 2011; Shang, 2017; Pinto-Llorente et al., 2017); (2) it considers a secondary education setting, also fundamental and rarely analyzed (Van der Kleij et al., 2015); (3) it examines the country of Uruguay, which has a strong digital learning presence starting from early grades due to its participation in the One Laptop Per Child Program; (4) it follows a quantitative analysis exploiting a large scale dataset (not only in terms of the number of students but also in teachers, classrooms and schools diversity) as opposed to most case studies which often include a very small number of students and classrooms (Ho, 2015; Shang, 2017); and (5) it uses both parametric and non-parametric models to account for possible non-linear effects of feedback characteristics on students' posting behavior and to avoid overly rigid statistical assumptions such as normally distributed residuals. In particular, our study draws inspiration from Miyamoto et al. (2015) who implemented a bootstrap for MOOC data analysis. Our paper provides another illustration of how this powerful method can be used in learning analytics and educational data mining research and how, in some ways, it is more flexible than parametric models.

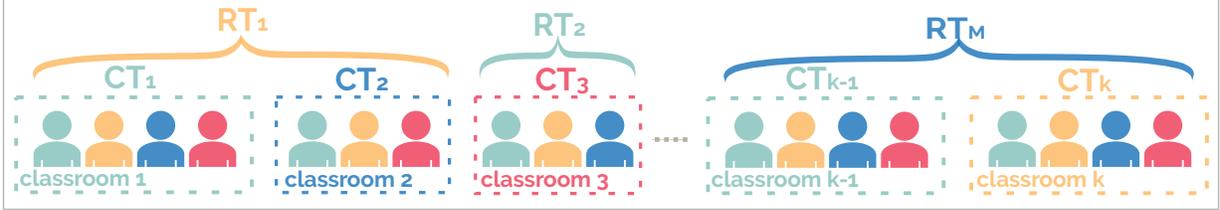


Figure 2: Nested structure of the dataset.

### 3 Dataset description

The dataset under consideration was originally collected by Aguerrebere et al. (2018) and includes all the comments (i.e., content, posting date, user ID) as well as administrative information (the CT, classroom, and RT for each student) for the secondary school classrooms (12-year-olds) that participated in the TDL program during school year 2017. In this work the dataset is extended to also include school year 2018. This includes a total of 27,627 comments exchanged between 1074 students, organized into 83 classrooms (in 49 public high schools located in 18 different states in the country), and 35 RTs. The dataset has a nested structure (see Figure 2): Students are organized into classrooms. Each RT, as well as each CT, can serve multiple classrooms. Figure 3 shows the histogram of the total comments posted by each student during their corresponding school year. The dataset has been de-identified to preserve each participant’s privacy and handled according to Uruguayan privacy-protection legislation. For privacy reasons, student grade data were not available. After talking with the TDL stakeholders and the program leaders, a set of features characterizing each comment was defined: *complexity*, *specificity*, *polarity* and *response delay*. Each feature represents a different aspect of how elaborate a comment is (*complexity*, *specificity*), its tone (*polarity*) and how long the student had to wait to receive feedback (*response delay*).

**Complexity** ( $c$ ) measures how elaborate a comment is, by adding its characters per word, words per sentence and total sentences:  $c = \frac{1}{4} \frac{\#char}{\#words} + \frac{1}{5} \frac{\#words}{\#sent} + \#sent$  (weights are included

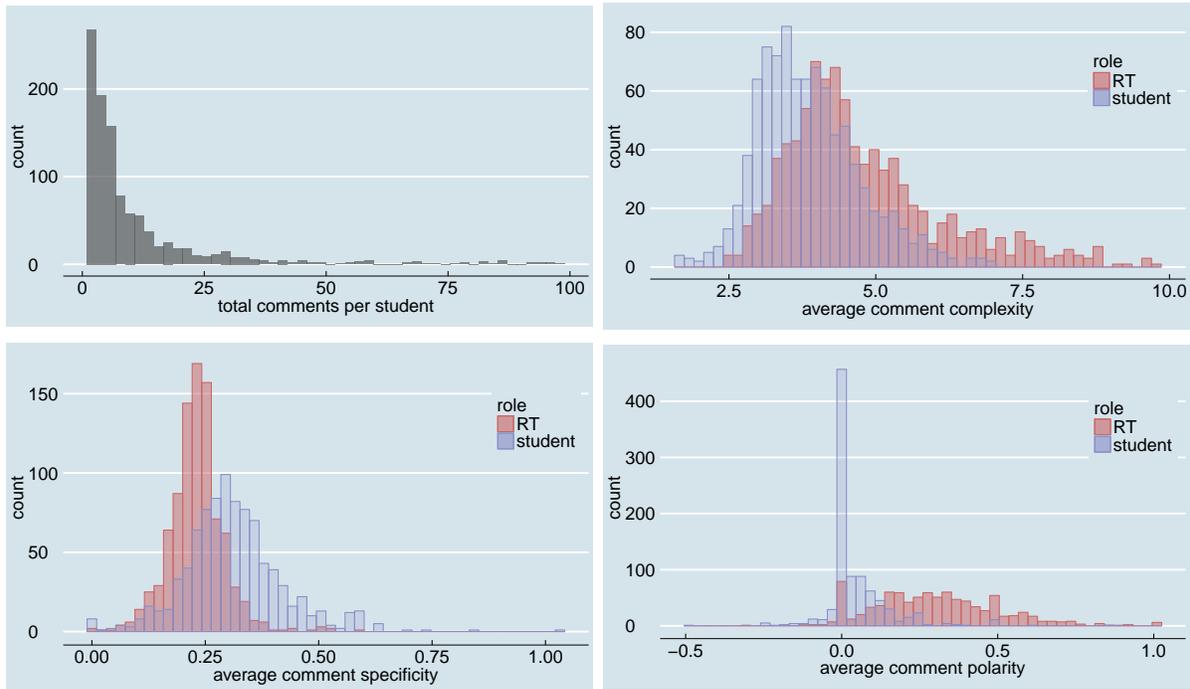


Figure 3: Histogram of the total comments posted by each student (top left). Histograms of the average complexity (top right), specificity (bottom left) and polarity (bottom right) of the RTs' feedback (red) and of the students' posts (violet).

to give similar relevance to all terms, 4 and 5 are the median characters per word and words per sentence respectively). Sentence boundaries were detected using the `sent.tokenize` of the `nltk` package. Examples of low, medium and high complexity comments:

*“Well done!”* ( $c = 2.4$ )

*“My favourite national park is Yellowstone.”* ( $c = 3.7$ )

*“Hi Alberto! This is an accurate description of the different continents, but can you try again? The activity is asking about different volcanic landforms! Can you please look at the encyclopedia and read the part about volcanic landforms to find the names of the three types of volcanic landforms? Here’s the link: [link].”* ( $c = 8.9$ )

This definition was motivated by, but is distinct from, the Automated Readability Index (ARI; Senter and Smith (1967)) that characterizes the readability of English texts. Like ARI, our measure includes terms for the average number of characters per word and the average number of words per sentence. We also included the number of sentences in our metric. We note that our definition is driven largely, though not entirely, by the length of the text. Many sentences in both students’ posts and teachers’ replies are often short, and thus the definition of complexity above does not amount simply to the total number of characters in the text. We decided not to use the ARI directly because it has low resolution and does not capture the variability within the TDL dataset. In particular, ARI is designed to approximate the United States grade level, and it is always an integer. Using ARI, almost all the comments in TDL would be mapped to a score of either 1 or 2 since the comments tend to be very basic in complexity.

**Specificity** ( $s$ ) measures how specific, on average, the words are in the comment. It combines how deep each word  $w_i$  appears in the WordNet (Miller, 1995) structure and how frequent the word is in the dataset:  $s = \frac{1}{W} \sum_{i=1}^W \frac{\text{depth}(w_i)}{Z} + \frac{1}{\text{freq}(w_i)}$ , where  $W$  is the total words in the comment and  $Z$  a normalizing factor equal to the maximum average comment complexity (Deshpande et al., 2010). Examples of comments with low and high specificity:

*“Very good! Do you have any cats?”* ( $s = 0.1$ )

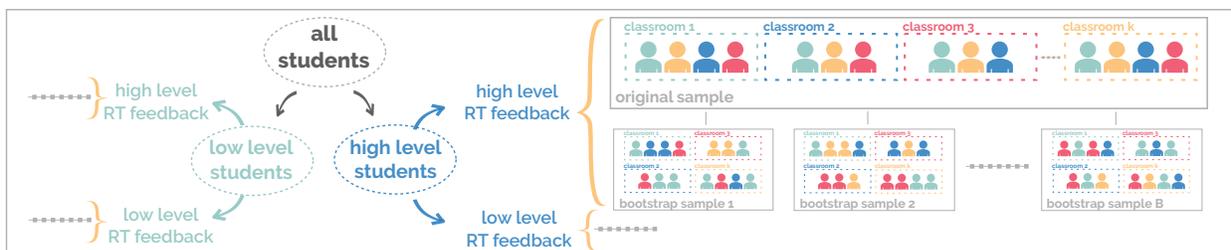


Figure 4: Bootstrap sample generation process used in Algorithm 1. Details are given for the case of high-level students who received high-level feedback, but the exact same process applies to the other three cases.

“*The skeleton of brontosaurus.*” ( $s = 1.2$ )

**Polarity** ( $p$ ) measures the tone of the comment (positive, negative) as the average of an index (-1 (negative) to +1 (positive))<sup>1</sup> assigned to each sentence based on the adjectives it contains (e.g., great, nice, awful). The accuracy of this system was reported as 74% (F-score) De Smedt (2013). Examples of positive and negative comments:

“*Great Carla! Awesome spelling!!*” ( $p = 1.0$ )

“*I would not like because they are dangerous.*” ( $p = -0.6$ )

**Response delay** ( $\tau$ ) is the timelapse between the student’s post and the RT’s response in days. Median  $\tau$  is 3 days.

Figure 3 shows the histograms of the average *complexity*, *specificity* and *polarity* of the RTs’ feedback and the students’ posts. The average specificity tends to be larger for students than for RTs; this is likely because students’ comments are responses to exercises that elicit very specific words such as “*skeleton of brontosaurus.*”, whereas the RTs’ comments often contain basic words, e.g., “*Great work!*”.

<sup>1</sup>The *sentiment* function of the PATTERN.EN Python module was used: <https://github.com/clips/pattern>.

## 4 Analyses

We first present a panel discussion we conducted with teachers to learn more about their intuition and experiences that might inform our analyses. Then we describe our data analysis procedures.

### 4.1 Group Discussions with Teachers

In order to learn more about the factors that could impact students' commenting behavior, especially within the particular learning context of the TDL program in Uruguay, we conducted a discussion session in December 2018 with about 30 TDL program coordinators and administrative staff, and about 10 remote teachers. The session was conducted in Spanish; a transcript of the participants' responses was translated into English. The format of the discussion was as follows:

1. We first asked the participants to respond to the question (Q1), "What do you think influences the students' posting behavior?"
2. Participants reflected on the question and wrote their answers on a sheet of paper.
3. We presented and explained the methodology and results of our statistical analyses (see sections below) and then asked the participants (Q2), "What do you think the causal mechanism behind our preliminary results might be?"
4. Participants reflected and wrote their answers on the paper.

All responses were anonymous. We presented our results *after* asking and obtaining answers to the first question so as not to bias their responses.

In total, 18 participants (both RTs and program coordinators who are well-familiar with the program's implementation) responded, though some did not respond to both questions.

A transcript of all the responses to both questions is given in Appendix A. We discuss the results of Q1 here and return to Q2 in the Discussion section at the end (Section 5). In terms of what factors were deemed to impact students' posting behavior, five main themes emerged:

- **Topics aligned with students' interests:** 12 participants (# 3,4,5,6,7,8,9,11,14,15,16,17) hypothesized that students post more comments when the topic of the exercise aligns with students' own personal interests.
- **Feedback characteristics:** 6 participants mentioned the characteristics of the RTs' feedback itself – such as latency, frequency, and consistency (#4,9,11,13) as well as the explicit posing of a question (#9,15).
- **Peers:** 4 participants (#4,7,8,11) mentioned the relationships with and encouragement from peers.
- **RT:** 3 participants (#2,11,12) listed the RTs' enthusiasm for and commitment to the TDL program; 5 more participants (#1,3,6,7,9) stated that the students' relationships with the RTs could be influential.
- **CT:** 6 participants (#1,3,6,7,10,11) mentioned the CTs' enthusiasm for and commitment to the program; 1 participant (#2) mentioned that students post more when the CT integrates it into the core curriculum (e.g., assigns homework to post a comment); and 1 participant mentioned the rapport between RT and CT (#12).

In the following sections of this paper, we explore quantitatively several of the factors mentioned by the TDL stakeholders that are listed above. For example, peer effects can be partly (though certainly not completely) captured by modeling the effect of the *classroom* on students' posting behavior. The effects of the RT and CT can be estimated by modeling them as random effects. Characteristics of the feedback itself (e.g., complexity, latency) are estimated using fixed effects.

The article is organized as follows. Sections 2 and 3 summarize the previous work and the

dataset under consideration. Sections 4.2 and 5 present the analysis with the corresponding discussion. Conclusions are summarized in Section 6. This paper extends the 6-page short paper with the same title that was published at the *Educational Data Mining* conference in 2019 Aguerrebere et al. (2019).

## 4.2 Statistical Modeling

We examine the effects of various feedback characteristics of RTs’ feedback on students’ posting behavior. One confound that must be considered is the student’s English proficiency level, as the answer to these questions may vary for more or less proficient students. We use two complementary approaches (Miyamoto et al., 2015): (1) multilevel linear regression that models the nested nature of the data; and (2) non-parametric bootstrap analysis. The latter is more complicated (e.g., requires a bin width parameter) but can model non-linear relationships and makes fewer assumptions (e.g., normality of residuals) than many parametric models.

## 4.3 How complex should the RT feedback be?

In this section we are interested in the question: Does RT feedback complexity (low vs. high – defined below) affect the total number of comments posted by the student? Note that we must consider the potential confound of the student’s English proficiency level, as the effect of RT feedback complexity may vary for more or less proficient students.

### 4.3.1 Non-parametric approach

To answer this question we first follow a non-parametric approach, using bootstrap to test the null hypothesis:

$$H_0 : E[T|S_c] = E[T|S_c, R_c], \quad (1)$$

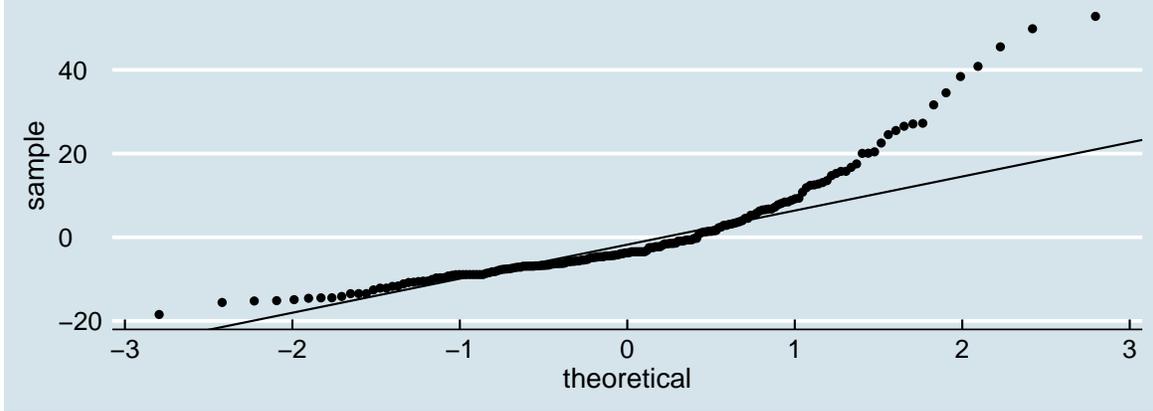


Figure 5: Q-Q plot of the residuals of the nested ANOVA fitted to the data for high-level students (similar results are obtained for low-level students).

versus  $E[T|S_c] \neq E[T|S_c, R_c]$ , where  $T$  is the total comments posted by the student,  $S_c$  is the student’s English proficiency level (low/high) and  $R_c$  is the complexity level of the feedback the student received from his RT (low/high). Hypothesis (1) tests whether the total comments posted by the student depend on the complexity level of the feedback she received from her RT, after conditioning on her English proficiency level. If the null hypothesis is rejected, then the complexity level of the feedback that the students receive affects their average engagement with the program, measured by the total comments they post and conditioned on their English proficiency level. It is important to note that we must reject the null hypothesis if  $E[T|S_c]$  differs statistically significantly from  $E[T|S_c, R_c]$  for *any* values of  $S_c$  and  $R_c$ . For this reason, in this section we examine the potential impact of  $R_c$  on  $T$  for each possible combination of  $(S_c, R_c)$ .

**Definition of high/low:** To estimate the student’s English proficiency level we use the average complexity of the comments posted by the student. Students with average complexity below (above) the median are classified as having low (high) proficiency respectively. Similarly, the complexity level of the RT’s feedback is low (high) when it is below (above) the median.

**Bootstrapping for equality of means:** How can we test whether  $P(T|S_c, R_c = low)$  and  $P(T|S_c, R_c = high)$  have the same mean? If the distribution  $P(T|S_c, R_c)$  were Gaussian for all values of  $S_c$  and  $R_c$ , then we could just use a  $t$ -test to compute the  $p$ -value (or a nested ANOVA to take into account the nested structure of the data). However, in our case the data are not Gaussian (see Figure 5). Fortunately, the bootstrap procedure proposed by Efron and Tibshirani (1994) provides a rigorous methodology: Let  $t_o$  be the normalized difference in means between  $P(T|S_c, R_c = low)$  and  $P(T|S_c, R_c = high)$  over the entire dataset:

$$t_o = (\mu_l - \mu_h) / \sqrt{\sigma_l^2/n_l + \sigma_h^2/n_h}. \quad (2)$$

where  $\mu$  represents the mean,  $\sigma$  the standard deviation, and  $n$  the number of samples in either the low-proficiency (subscript  $l$ ) or high-proficiency (subscript  $h$ ) data subset. Recall that the mathematical definition of  $p$ -value is the probability, *conditional* on the null hypothesis  $H_0$ , of observing a test statistic at least as large as the observed statistic (in our case,  $t_o$ ). By sampling *with replacement* from our original dataset, we can *simulate* multiple data samples. We subselect the data for which  $R_c = low$  and the data for which  $R_c = high$  and then resample each of them to generate multiple bootstrap samples. To enforce the null hypothesis (i.e., equal means), we *set* the means of the two samples to be equal to the mean of the combined sample. We then compute the normalized difference in means between the two subsets in the bootstrapped sample. Over all  $B$  bootstrap iterations ( $B = 10000$ ), we finally compute the fraction in which the normalized difference in means is at least as large as the observed statistic. Algorithm 1 presents the pseudo-code of the bootstrap method used to test Hypothesis (1). Note that this method is more suitable to our case than a classical permutation test, which tests the hypothesis of equal *distributions* (not just the first moment) and thus assuming equal means *and* variances. As there is no compelling reason to assume equal variances in our case, we opt to test the equality of means only.

---

**Algorithm 1: Bootstrapping Hypothesis Testing**

---

```
1 for  $i$  in  $[low, high]$  do /* Student complexity level */
2   for  $b = 1, \dots, B$  do /* Bootstrap samples */
3     for  $j$  in  $[low, high]$  do /* RT complexity level */
4       Generate bootstrap sample: sample with replacement the students within each
5         classroom, whose  $S_c = i$  and  $R_c = j$  (see note2).
6       Compute  $\mu$  and  $\sigma^2$  (average and variance of  $T$  for the mean-adjusted bootstrap
7         sample, see Section 4.3.1).
8     end
9     Compute test statistic:  $t_b = (\mu_l - \mu_h) / \sqrt{\sigma_l^2/n_l + \sigma_h^2/n_h}$ , with  $n_l, n_h$  the bootstrap sample
10    sizes and subindex  $l$  and  $h$  corresponding to low and high RT complexity levels
11    respectively.
12  end
13  Compute  $p$  as  $length(abs(t_b) \geq t_o) / B$ , with  $t_o$  the observed test statistic value.
14 end
```

---

**Results:** For high-level students, more *basic* feedback is associated with more posting. Students who received high-level feedback posted on average 22% fewer comments than those who received low-level feedback (9.3 versus 12 total comments,  $p = 0.006$ ). Examples:

**Student A:** *My favorite food is hamburguer.*

**RT (low level):** *Nice Lucia! What do you eat in your hamburger?*

**Student B:** *They are in Spain in Barcelona.*

**RT (high level):** *Hello Pablo! Yes they are in Spain. Very good. Here is a link in case you would like to know a bit more about Spain and their culture. In Uruguay there are a lot of people who have Spanish origins. It is very evident in the food :) I have never been to Spain. Would you like to go to Spain?*

A possible explanation for this behavior is that even the *high-level* students have weak English proficiency and might feel overwhelmed by feedback that is too complex. No statistically significant differences are observed for low-level students (8.5 versus 7.6 total comments,  $p = 0.14$ ).

### 4.3.2 Controlling for the RT

Another possible confound is the effect of the RT him/herself, as more motivated RTs may give better feedback that leads to more posts by their students. To take this into account,

---

<sup>2</sup>A minimum classroom size is imposed to ensure bootstrap samples of at least 5 students.

null Hypothesis (1) is reformulated as:

$$H_0 : E[T|S_c, RT] = E[T|S_c, R_c, RT], \quad (3)$$

versus  $E[T|S_c, RT] \neq E[T|S_c, R_c, RT]$ , where  $RT$  is the ID of the student’s remote teacher. If Hypothesis (3) is rejected, the complexity level of the feedback *given by a particular RT* affects the total comments posted by the student. A variation of the bootstrapping algorithm is used to test Hypothesis (3) where, instead of defining the  $[low, high]$  RT feedback complexity levels globally from all samples, independent thresholds are defined for each RT (Algorithm 2).

---

**Algorithm 2:** Bootstrapping Conditioning on RT

---

```

1 for  $i$  in  $[low, high]$  do                                     /* Student complexity level */
2   for  $b = 1, \dots, B$  do                                       /* Bootstrap samples */
3     for  $r = RT_1, \dots, RT_M$  do                                   /* for each RT */
4       Define  $[low_r, high_r]$  complexity levels
5       for  $j$  in  $[low_r, high_r]$  do                               /* RT complexity level */
6         Generate bootstrap sample: sample with replacement the students within each
           classroom, whose  $RT = r, S_c = i, R_c = j$ .
7         Compute  $\mu$  and  $\sigma^2$  (average and variance of  $T$  for the mean-adjusted bootstrap
           sample, see Section 4.3.1).
8       end
9       Compute test statistic:  $t_b^r = (\mu_l - \mu_h) / \sqrt{\sigma_l^2/n_l + \sigma_h^2/n_h}$ , with  $n_l, n_h$  the bootstrap
           sample sizes and subindex  $l$  and  $h$  corresponding to  $low_r$  and  $high_r$  RT complexity
           levels respectively.
10      end
11      Compute average test statistic across RTs:  $t_b = \frac{1}{M} \sum_{r=1}^M t_b^r$ 
12    end
13    Compute  $p$  as  $length(abs(t_b) \geq t_o) / B$ , with  $t_o$  the observed test statistic value.
14 end

```

---

**Results:** The result previously obtained remains valid even when conditioning on the RT (i.e., resampling within each RT), meaning that different feedback levels given by the RT to his students are associated with different total posts (8.9 versus 12.5 total comments,  $p = 0.01$ ). No statistically significant differences are observed for low-level students.

Even if controlling for the RT makes the result more solid from a statistical point of view,

interpreting the results becomes more difficult, as the meaning of *low* and *high* complexity feedback changes from RT to RT. Because a fundamental goal is to translate these results into useful information for the teachers, this approach has the disadvantage that an *absolute* complexity level reference cannot be given to them as reference of what *low* and *high* means, and how to position the feedback they give with respect to that.

### 4.3.3 Parametric Approach

A parametric approach is conducted to complement the results obtained by the non-parametric analysis. The parametric model allows for the inclusion of additional predictors to avoid possible confounds; it enables estimating the linear “dose response” of different feedback characteristics; and it does not require binning of RT complexity (though low/high bins are still used for student complexity). The price we pay for these capabilities is that they make more rigid assumptions of normally distributed residuals and linearity in the predictors. We restrict the analysis to classrooms with  $\geq 10$  students so we can compare with the non-parametric approach.

In order to take into account the nested structure of the data, a multilevel modeling approach is employed where the classroom and RT effects on the student’s activity are modeled as nested random effects. Note that the *classroom* of a student is distinct from the *classroom teacher* (CT) because each CT can potentially teach multiple classrooms. Hence, the classroom ID can model peer effects more directly than it models the effectiveness of a particular teacher. (We return to the issue of modeling the CT as a random effect in Section 4.6).

We model student  $i$ ’s total posts as a negative binomial random variable, to account for the fact that it is count data with overdispersion, with expected value  $\mu_i$  given by:

$$\log(\mu_i) = \beta + \gamma_0 c_i + \gamma_1 y_i + \gamma_2 p_i + \gamma_3 s_i + \gamma_4 \tau_i + K + R, \quad (4)$$

	stud. level	$\beta$		$\hat{\gamma}_0$		$\hat{\gamma}_1$		$\hat{\gamma}_2$		$\hat{\gamma}_3$		$\hat{\gamma}_4$
Model 4	low	1.12	**	0.03		0.30	.	-0.03		0.76		0.01
	high	2.50	***	<b>-0.11</b>	*	0.26		-0.02		-0.02		0.04
Model 5	low	1.63	***	<b>-0.14</b>	**	0.25		-0.32		0.67		-0.01
	high	2.24	***	<b>-0.16</b>	**	0.22		-0.06		0.12		-0.002

Table 1: Effects of RT feedback on students’ total comments for Models 3 and 4, in log scale. Signif. codes: 0 (\*\*\*) 0.001 (\*\*) 0.01 (\*) 0.05 (.) 0.1

(capital letters denote random variables and lower-case denote fixed values).  $\beta$  is the baseline total comments.  $c_i$ ,  $p_i$ ,  $s_i$  and  $\tau_i$  are the average complexity, polarity, word specificity and response delay of the feedback comments student  $i$  received from his RT, respectively.  $y_i$  is the school year. The fixed effects  $\gamma_0, \dots, \gamma_4$  represent the effects of the corresponding covariates on the total comments. The nested random effect classroom-RT is represented by the random variables  $K$  and  $R$ , assumed to follow zero-mean Gaussian distributions with standard deviations  $\sigma_K$  and  $\sigma_R$ . All the parametric models were fit using the R *lme4* package (Bates et al., 2015).

**Results:** Table 1 (Model 4) shows the computed effects for all the covariates. The parametric analysis, which includes other possible confounds, confirms the same tendency observed with the non-parametric approach: a negative statistically significant effect of the RT feedback complexity level is observed for high-level students ( $\exp(-0.11) = 0.9$ , i.e., 10% less total posts per unit increase in RT average complexity) and no effect is observed for low-level students. To compare this result to the one obtained by the non-parametric approach, we compute the equivalent per unit decrease in the non-parametric case (computed as the total decrease divided by the difference between the average RT complexity in the two compared levels, 3.8 and 5.9) which equals 10%.

## 4.4 Should RTs feedback complexity be close to that of their students?

Rather than the *absolute* complexity, we can also consider the *relative* complexity of the RTs' feedback compared to the complexity of students' comments. Put another way: should the feedback complexity be somehow *adapted* to the student? To answer this question we propose to model the total comments posted by the student as a function of the *distance* between the average complexity of the student's comments and the average complexity of the feedback the student received from his RT.

### 4.4.1 Parametric approach

We model each student  $i$ 's total posts as a negative binomial random variable with expected value  $\mu_i$  given by:

$$\log(\mu_i) = \beta + \gamma_0|c_i - c_{s_i} - \alpha| + \gamma_1y_i + \gamma_2p_i + \gamma_3s_i + \gamma_4\tau_i + K + R, \quad (5)$$

The fixed effect  $\gamma_0$  represents the effect of the absolute value of the difference between the student's ( $c_{s_i}$ ) and the RT's ( $c_i$ ) average comments complexity, where  $\alpha$  is introduced as an offset to account for the fact that the feedback may need to be close to that of the student but not necessarily equal. See Model 4 for the definition of the rest of the variables.

**Setting  $\alpha$ :** Model 5 is fitted for different values of  $\alpha$  and the one corresponding to the largest log-likelihood is selected. For low student levels the maximum log-likelihood is obtained at  $\alpha = 0.25$ , whereas for high student level it is at  $\alpha = -0.79$ . Hence, this analysis suggests that even if for both low and high-level students feedback complexity should be close to the student level, low-level students benefit from feedback slightly more complex than theirs whereas high-level students benefit from feedback slightly below theirs. Recall that low-level students post very basic comments, whereas those of high-level students tend

to be more elaborate but remain still simple.

**Results:** Table 1 (Model 5) shows the computed fixed effects for the different covariates. The distance between the average complexity of the student's comments and the average complexity of the feedback the student received from his RT has a negative stat. sig. effect on the total comments posted by the student. There is a 13% ( $p = 0.008$ ) and 15% ( $p = 0.003$ ) decrease in total comments with one unit distance increase, for the low and high-level students respectively. We present in the following a series of examples in order to help the reader gain insight into what small and large student-RT complexity distance mean in practice.

### **Large distance:**

**Student A:** *I would like to defile.*

**RT:** *Nice try Marcela, but I do not understand what you mean. There are two models in the above photo, can you tell me which model is from Brazil and which model is from the USA? or - can you tell me, who is your favourite model? My favourite model is Kate Moss.*

**Student B:** *I like to watch suspense movie.*

**RT:** *Gabriela! I don't know any suspense movies. Do you have a favourite?? The only move I can think of that is suspense is this one, Splice, its very weird and kind of creepy. But I really liked it!! Tell me what you think :) [link]*

The RTs' responses are complex compared to what the student wrote.

### **Small distance:**

**Student A:** *My favorite sport is football.*

**RT:** *Very good! what is your favorite football team?*

**Student B:** *I believe that the aliens live in Jupiter.*

**RT:** *Why Jupiter and not Mars?*

The Student's and RT's comments have comparable complexity levels and, unlike the *large distance* examples, the interactions seem closer to what a regular online chat could be like

for students with basic English proficiency.

#### 4.4.2 Non-parametric approach

A non-parametric approach is conducted to complement the results obtained by the parametric analysis. For this purpose, bootstrapping is used to test the null hypothesis:

$$H_0 : E[T|S_c] = E[T|S_c, D], \quad (6)$$

versus  $E[T|S_c] \neq E[T|S_c, D]$ , where  $D$  is the distance between the student's and the RT's average comments complexity as defined in Model 5. The bootstrapping algorithm introduced in Section 4.3.1 is used, with a for loop on small and large  $D$  (below and above the median distance) instead of RT complexity, and  $\alpha$  is set to the values obtained in Section 4.4.1.

**Results:** A negative statistically significant effect is observed for  $D$  on the total comments posted by the student, both for low and high-level students. Students who received feedback less adapted to their level posted 15% and 34% less comments than those who received more adapted feedback, for low (6.6 versus 7.6 total comments,  $p = 0.007$ ) and high (9.2 versus 12.3 total comments,  $p < 0.001$ ) level students respectively. To compare these results to those obtained by the parametric approach we compute the equivalent per unit decrease (computed as the total decrease divided by the difference between the average  $D$  in the two compared levels) which equals 9% both for low and high student level.

### 4.5 Engaging students in conversation

The program aims at motivating the students to interact with others in English. Therefore, we are interested not only in their total posts but also in the probability of engaging them in a conversation, which we operationalize as the student responding to the teacher's response to the student's original message. Recall that the students are always the ones that start

the conversation threads by posting the first comment about a given topic. Then, the RT responds, and the student may or may not continue the conversation in the same thread. The non-parametric bootstrap conducted previously is not suited to this case because we would need to compute the probability that a given *student* engages in conversation. However, with our definition, engagement in a conversation is not tied to a student but rather to each thread. The probability could be approximated by the ratio between the total threads where the student posted at least a second comment and the total threads initiated by the student. Setting aside that this is not exactly the variable of interest (i.e., the probability that a student posts at least a second comment in a given *thread*), this approximation will be noisy and hard to interpret as, for instance, it would take the same value (1.0) for a student who continued all his 100 conversations and a student that started and continued only 1 conversation. Therefore, we opt to conduct a multilevel regression approach which is well adapted to this section’s goal.

Following the same rationale as Section 4.3.3, we use a multilevel logistic regression model to explore this question. Let  $(Y_j)_{j=1,\dots,N}$  be Bernoulli variables with  $P(Y_j = 1|\eta_j) = \exp(\eta_j)/(1 + \exp(\eta_j))$  with:

$$\begin{aligned} \eta_j = & \beta + \gamma_0 c_j + \gamma_1 \log(\tau_j) + \gamma_2 p_j + \gamma_3 s_j + \gamma_4 y_j + \\ & \gamma_5 q_j + \gamma_6 l_j + \gamma_7 e_j + S + K + R, \end{aligned} \tag{7}$$

$Y_j = 1$  if the student posts a second comment in conversation  $j$  and 0 otherwise.  $\beta$  is the baseline.  $c_j$ ,  $p_j$  and  $s_j$  are the complexity, polarity and specificity of the RT’s response to the first comment posted by the student who initiated conversation  $j$ .  $q_j$ ,  $e_j$  and  $l_j$  are boolean variables taking value 1 if the RT’s response asked the student a question, included an emoticon or shared a link respectively.  $\tau_j$  is the timelapse between the moment the student started conversation  $j$  and the RT replied.  $y_j$  is the school year.  $\gamma_0, \dots, \gamma_7$  represent the fixed

student level	$\beta$	$\hat{\gamma}_0$	$\hat{\gamma}_1$	$\hat{\gamma}_2$	$\hat{\gamma}_3$	$\hat{\gamma}_4$	$\hat{\gamma}_5$	$\hat{\gamma}_6$	$\hat{\gamma}_7$	
low	-7.11 ***	0.08	-0.15 ***	0.36	-2.20 **	0.24	2.94 ***	-0.53	0.11	
high	-7.46 ***	0.35 ***	0.10 **	-0.47 *	-2.66 **	0.62	2.83 ***	-0.77 .	-0.39	
low (question asked)	-3.78 ***	0.00	-0.16 ***	0.83 **	-3.05 **	0.25	-	-	-0.30	0.13
high (question asked)	-4.27 ***	0.22 **	0.08 .	0.59 .	-2.76 **	0.62	-	-	-0.63	-0.84 .

Table 2: Effects of RT feedback characteristics on the probability of engaging the student in conversation, in logarithmic scale. Significance codes: 0 (\*\*\*) 0.001 (\*\*) 0.01 (\*) 0.05 (.) 0.1

effects of the corresponding covariates. The nested random effect student-classroom-RT is represented by the random variables  $S$ ,  $K$  and  $R$ , assumed to follow zero-mean Gaussian distributions with standard deviations  $\sigma_S$ ,  $\sigma_K$  and  $\sigma_R$ .  $N$  is the total number of conversation threads and there may be several threads per student.

**Results:** Table 2 shows the estimated covariate effects. By far, and maybe not surprisingly, the fact that the RT asks the student a question has the largest stat. sig. positive effect on the probability of getting the student to continue the conversation. When a question is asked, assuming the rest of the covariates remain fixed, the odds ratio for students of the same classroom is  $\exp(2.94) = 19.0$  and  $\exp(2.86) = 17.5$ , for low and high-level students respectively.

Complexity and polarity are negatively correlated (Spearman correlation -0.5 for high and -0.23 for low-level students respectively), as very positive comments such as “*Great work!*” tend to be quite basic in terms of complexity. This may explain the opposite signs of the stat. sig. effects observed for complexity and polarity for high-level students. As more elaborate comments often include questions, the positive effect of complexity suggests that more elaborate comments increase the probability of engaging a student in a conversation. On the contrary, the negative stat. sig. effect of polarity is likely due to the fact that very positive comments such as “*Great work!*” tend to be quite basic in terms of complexity. For high-level students a larger delay is associated with more responses, as after a one-week delay the odds ratio is 1.2 (the rest of the covariates and random effects remaining constant).

This is likely because for high-level students RT response delay is positively correlated with complexity (Spearman 0.1) and negatively correlated with comments polarity (Spearman -0.12): writing more elaborate comments take longer. This is not observed for low-level students (Spearman 0.02 for both complexity and polarity), for whom it seems important to reply as soon as possible, as after a one-week delay the odds ratio is 0.7. Finally, for both high and low-level students, results suggest that using less specific words is associated with higher probability of engagement in the conversation.

**What if the RT asked a question?** Because the probability of getting the students to continue the conversation highly increases when they are asked a question, we would like to know if the RT feedback characteristics have the same effect *given a question was asked*. For this purpose, we repeat the analysis conditioning on the fact that a question was asked by the RT, and fit Model 7 in the dataset subsample that verifies  $q_j = 1$ . The third and fourth rows of Table 2 show the obtained results, which are consistent with those obtained with the complete dataset. The only difference is that polarity now has a positive effect for both low and high-level students. This may be explained by the fact that restricting the analysis to RT comments which include a question leaves out the very basic RT posts such as “*Thanks!*” or “*Nice work!*”. Hence, assuming a minimum complexity level given by the fact that at least a question is asked, more positive comments are associated with an increased probability of engaging the student in the conversation.

## 4.6 Effect of the CT

In the particular context of the TDL program, the RT and CT work together to help the students to learn English: the CTs know their pupils personally and interact with them twice a week in the same physical space. They are responsible for implementing the EFL pedagogy (explaining concepts, answering questions, etc.) and evaluating students. They are

required to be fluent in English but are rarely native speakers. In contrast, the RTs spend much less time with the students than the CTs and are separated physically, which can alter the interpersonal dynamics. They are required to be fluent in English, and are usually native speakers. RTs and students communicate once a week through videoconference in a conversation class, and online using discussion forums. Within this “learning triad”, the question arises of what impact the CTs have on students’ learning. One possible impact is that some CTs are more enthusiastic about the TDL program than others and may encourage their students to participate online in the exercises (including commenting). While we do not have a direct measure of learning, we can at least explore the impact on students’ commenting activity. To this end, we devised a parametric model that included both the CT and RT as random effects. In contrast to the model in Section 4.3.3, in which the *classroom* was captured as a random effect  $K$ , here the classroom *teacher* is captured as a random effect  $C$ ; the distinction is that each CT may teach in multiple classrooms.

$$\log(\mu) = \beta + C + R, \tag{8}$$

By fitting this model to the TDL data and then comparing it to a similar model *without* the  $C$  term,

$$\log(\mu) = \beta + R, \tag{9}$$

we can assess whether the CT random effect  $C$  is statistically significant, and also get a sense of how large this impact is compared to the overall variance of the data. As in the other experiments, we fit this model separately for high- and low-proficiency students. To compute  $p$ -values for the significance of the estimated variance  $\sigma_C^2$  of the CT random effect  $C$ , we used a parametric bootstrap approach (Algorithm 3). The essence is that we compute a null distribution for the difference in deviance (a measure of goodness-of-fit), between Models 8 and 9 w.r.t. the actual dataset, by assuming that the “true” generative process for

the number of comments does not include the CT as a random effect (i.e., Model 9).

---

**Algorithm 3:** Parametric Bootstrapping for CT Random Effect

---

```

1 for  $i$  in  $[low, high]$  do                                     /* Student complexity level */
2   Fit Model 9 to estimate variance  $\sigma_R^2$  of  $R$ .
3   for  $b = 1, \dots, B$  do                                     /* Bootstrap samples */
4     Using Model 9, sample values for each RT using the variance estimated in step (1); then,
       for each data point in the actual dataset, sample the number of comments from the
       negative binomial distribution with expected value  $\mu$ . Call this dataset  $b$ .
5     Fit Model 9 using bootstrap dataset  $b$ .
6     Fit Model 8 using bootstrap dataset  $b$  (where the CT of each data point is obtained from
       the actual dataset).
7     Compute the deviance for each model w.r.t. dataset  $b$ , and then compute the difference  $\Delta_b$ 
       in deviances.
8   end
9   Compute  $p$  as  $length(\Delta_B \geq \delta)/B$ , with  $\delta$  the observed difference in deviance between Models 9
       and 8 on the actual dataset.
10 end

```

---

**Results:** For students with low English proficiency, the variance  $\sigma_C^2$  of the CT random effect was estimated as 0.27; this variance is statistically significant ( $p = 0.006$ ). For students with high English proficiency,  $\sigma_C^2 = 0.38$  and was also significant ( $p = 0.007$ ). Note that this variance impacts the *logarithm* of the expected value  $\mu$  of the negative binomial distribution, not the expected number of comments itself. Hence, to put these results into perspective, we can compute the expected number of comments for a student with a “low TDL activity” CT as well as the expected number of comments for a student with a “high TDL activity” CT, where we operationalize “low/high TDL activity” as the 5th and 95th percentiles of the Gaussian distribution for  $C$ , respectively. We then compare this range to the overall 5%/95% range of number of comments in the observed dataset. This gives a sense of how much of the total variance is explained by just the CT. For high-proficiency English learners, the estimated range due to the CT is [2.4, 28.9], whereas the 5%/95% range for the actual dataset is [0,31]. For low-proficiency learners, the estimated range due to the CT is [1.7, 13.3], whereas the 5%/95% range for the actual dataset is [0,20]. Hence, we see that the impact on the expected number of comments due to just the CT can be quite substantial.

## 5 Discussion

How might the insights gained from this study inform teaching methods and influence feedback models for education technology platforms? First, let us recall that no *causal* inference can be made directly from this study as, even if the RT and CT effects and other important confounds were taken into account, there may always be other factors influencing the observed behavior. It did generate, however, interesting and plausible hypotheses to be evaluated through an experimental design.

That said, the obtained results suggest that the RTs should pay special attention to their students' English proficiency level, by observing how complex their comments are, and try to adapt their feedback accordingly. Providing feedback that is too complex or too basic seems to be counterproductive for the students. This is not only against the original goals of the program itself (to motivate the students to participate and post as many comments as possible) but it may also be a waste of resources as writing too elaborate responses take the RTs much more time and effort. Moreover, this may result in longer response delays and less motivated RTs. We remind the reader that our definition of complexity is driven primarily, but not entirely, by the length of a discussion forum post. A limitation of the present study is that its results are largely tied to the specific measure we used. Future work should explore alternative definitions of complexity (e.g., Biber (1992); Graesser et al. (2014)) that consider the grammatical structure of the texts as well.

During the meeting with some TDL coordinators and RTs, we presented some of the preliminary results of this study and asked the participants: *What do you think the causal mechanism behind our preliminary results might be?* (Q2). The participants suggested several reasons why students may post more comments when their RTs' feedback complexity is similar to their own: (1) Similar complexity means that students have a higher chance of understanding what the RT meant, compared to if the RT complexity is much higher

(#3,4,8,9). (2) Similar complexity may prevent the student from feeling overwhelmed or frustrated, while still providing an engaging challenge (#2,4). (3) RTs' ability to *adapt* might be correlated with their ability to *teach*, i.e., perhaps it was not the adaptation in complexity per se that mattered but simply the RTs' pedagogical effectiveness in general (#1).

Regarding how to engage the students in a conversation, results suggest (unsurprisingly) that the best way is to pose a direct question. Aligned with the previous results, a balance is needed between too elaborate and too basic comments. The comments should be complex enough to include for instance a question (i.e., very positive posts such as “*Well done!*” or “*Great work!*” won't be enough) yet keep their simplicity in terms of using mostly basic words (low specificity). For the students with the lowest English proficiency level, it seems important to reply as quickly as possible in order to keep them engaged.

## 6 Summary and Conclusions

We conducted an analysis of 27,627 comments, exchanged between 1074 high school students and 35 RTs over 2 years, to study the effect that different RT's feedback characteristics have on the students' posting behavior in an online EFL learning environment. While our analysis was observational, we controlled statistically for important potential confounds such as (1) the classroom teacher, (2) the remote teacher and (3) the students' English proficiency level. Through a combination of both parametric and non-parametric models, we also avoided making rigid assumptions about linearity of the fixed effects and Gaussianity of residuals. The research questions, as well as the features defined for the characterization of the comments, were discussed and validated with the stakeholders and the leaders of the program in order to take advantage of their wide experience on the topic. Our results suggest that:

(1) Teachers should observe the complexity of their students' comments and adapt the com-

plexity of their feedback accordingly. Students who receive feedback that is too complex or too basic for their level post 13% ( $p = 0.008$ ) and 15% ( $p = 0.003$ ) fewer comments than those who receive adapted feedback, for low and high-level students respectively.

(2) According to some RTs who were consulted about the potential causes of the observed behavior, the students may be more motivated when the language of the RT is accessible to them because they understand it, they learn from it and are challenged by it, without this turning into frustration.

(3) The best way to engage the students in a conversation is to pose a question (this increases the odds by 19 and 17.5 for low and high-level students respectively). The comments should be complex enough to include a question (i.e., “*Great work!*” won’t be enough) yet remain simple in terms words specificity. Also, for low-level students, it is important to respond as quickly as possible (after a one-week delay the odds ratio is 0.7).

(4) Finally, in the context of distance learning programs such as TDL, the classroom teacher (CT) him/herself can have a significant effect on students’ behavior.

Even if no strong causal inferences can be made, this study generated enlightening insights which have potential implications for the countless online platforms offering foreign language learning services, in which it is crucial to give the best possible learning experience while judiciously allocating resources (e.g. teachers’ time).

As future work, we would like to explore other feedback characteristics that may also have an effect on the students’ motivation (not just participation) within the program. It would also be interesting to analyze the evolution of students’ language proficiency over time and to assess whether this could be a sign of actual learning, as well as to study how RT feedback complexity should increase over time. Furthermore, an interesting question would be: Is it possible to automatically guide the RTs on *how* and *when* should the feedback be given to each student (based on the historical data of the student’s English proficiency level)?

**Acknowledgments:** We thank the TDL team and RTs for their participation and helpful feedback. J. Whitehill was supported by the National Science Foundation under Grant No. #1822830.

## A Appendix: Group Discussion with Teachers

Id	“What do you think influences the student’s posting behavior?”	“What do you think the causal mechanism behind our preliminary results might be?”
1	<ul style="list-style-type: none"> <li>• Student’s relationship with RT.</li> <li>• The RT’s encouragement of the students.</li> <li>• The CT’s encouragement of the students.</li> <li>• Desire to win a prize.</li> </ul>	<ul style="list-style-type: none"> <li>• I think the RTs adapt the feedback they give the students based on what the students post. Those RTs who are the best at encouraging the students are also best at adapting their feedback. Hence, I think the third variable of RT sensitivity/engagement is driving both the number of comments the students post and the feedback students receive.</li> </ul>
2	<ul style="list-style-type: none"> <li>• RTs are more comfortable with TDL; therefore, they promote it more.</li> <li>• CTs are including it as a part of the curricula.</li> <li>• RTs are giving more communicational feedback.</li> </ul>	<ul style="list-style-type: none"> <li>• It’s easier for students to connect when they feel they’re being understood (similar level).</li> <li>• It’s motivating to have some challenge (give more complex feedback), but not so big challenge that it turns frustrating.</li> <li>• If students are copying answers, we can assume that they know (...) but it’s better to be in a safe area.</li> </ul>
3	<ul style="list-style-type: none"> <li>• RT and CT encouragement.</li> <li>• Engaging topics.</li> <li>• Motivation.</li> <li>• Appealing visuals.</li> <li>• Training / Coaching.</li> </ul>	<ul style="list-style-type: none"> <li>• If students understand their feedback, they would respond better and more frequently.</li> </ul>
4	<ul style="list-style-type: none"> <li>• Topic / Content.</li> <li>• RT’s response / feedback.</li> <li>• Take-up/engagement by peers.</li> </ul>	<ul style="list-style-type: none"> <li>• Students are motivated when the language of the interlocutor is close to that of their own. 1) They understand it, 2) they learn/are challenged (comprehensible input).</li> </ul>

5	<ul style="list-style-type: none"> <li>• Their interest (interesting topics)</li> </ul>	<ul style="list-style-type: none"> <li>• Maybe RTs understand better their students.</li> <li>• 2018 have better level of English as they were more exposed to English lessons at school.</li> </ul>
6	<ul style="list-style-type: none"> <li>• RT's feedback and encouragement.</li> <li>• CT's support and commitment.</li> <li>• Students' interests.</li> </ul>	<ul style="list-style-type: none"> <li>• Complexity of the tasks.</li> </ul>
7	<ul style="list-style-type: none"> <li>• Their interests.</li> <li>• Motivation from the RTs and CTs and from their classmates.</li> </ul>	<ul style="list-style-type: none"> <li>• RTs have knowledge of the student's real English level and the context so it can adjust the feedback to each student's situation.</li> </ul>
8	<ul style="list-style-type: none"> <li>• RT asking follow-up questions.</li> <li>• Peers' feedback such as likes.</li> <li>• Engaging topics</li> </ul>	<ul style="list-style-type: none"> <li>• Input must be suitably graded for students to understand and respond to it.</li> </ul>
9	<ul style="list-style-type: none"> <li>• Relationship with RT.</li> <li>• Response time.</li> <li>• Topic.</li> <li>• Asking questions.</li> </ul>	<ul style="list-style-type: none"> <li>• Students are more likely to interact/post if the level of language is accessible to them.</li> </ul>
10	<ul style="list-style-type: none"> <li>• It depends a lot on how much the CT influences on them. The winner we had was because CT total take TDL as part of the program not taken as homework.</li> <li>• I also believe that RTs' responses are crucial as well.</li> </ul>	<ul style="list-style-type: none"> <li>• I agree that the balance in the RT's comments can be crucial.</li> </ul>
11	<ul style="list-style-type: none"> <li>• CT and RT motivation with the program.</li> <li>• Topics.</li> <li>• Comments and how fast are those comments posted by RT and partners.</li> <li>• Level of English.</li> </ul>	<ul style="list-style-type: none"> <li>• RTs and CTs have more experience in project now. They better know the students.</li> <li>• Both RTs and CTs have devised ways to better motivate students.</li> </ul>
12	<ul style="list-style-type: none"> <li>• RT's commitment to the program.</li> <li>• RT-CT rapport both in class and in the TDL space.</li> <li>• CT-RT rapport+coordination touching upon topics in class.</li> </ul>	<ul style="list-style-type: none"> <li>• Better RT-CT training</li> <li>• experience.</li> <li>• Knowledge of program + topics.</li> </ul>

13	<ul style="list-style-type: none"> <li>• Type of task</li> <li>• RTs feedback.</li> <li>• Frequency (consistency).</li> </ul>	<ul style="list-style-type: none"> <li>• Level of difficulty.</li> </ul>
14	<ul style="list-style-type: none"> <li>• Motivation.</li> <li>• Interests.</li> <li>• Engagement.</li> </ul>	<ul style="list-style-type: none"> <li>• The change in strategy.</li> </ul>
15	<ul style="list-style-type: none"> <li>• Motivating tasks (a reason for writing).</li> <li>• Getting responses from RTs, e.g., follow-up questions.</li> </ul>	
16	<ul style="list-style-type: none"> <li>• RT's responses.</li> <li>• Interesting topics.</li> </ul>	
17	<ul style="list-style-type: none"> <li>• Interesting topics for students.</li> <li>• RTs motivate students to participate.</li> </ul>	
18		<ul style="list-style-type: none"> <li>• RTs better trained on how and what kind of feedback to give.</li> <li>• More awareness of how to shift feedback to students' level.</li> </ul>

## References

- Aguerreberre, C., Bulger, M., Cobo, C., García, S., Kaplan, G., and Whitehill, J. (2019). How should online teachers of english as a foreign language (efl) write feedback to students? In *Proceedings of the Educational Data Mining Conference*.
- Aguerreberre, C., Cabeza, S. G., Kaplan, G., Marconi, C., Cobo, C., and Bulger, M. (2018). Exploring feedback interactions in online learning environments for secondary education. In *Proc. of the 1st Latin American Workshop on Learning Analytics*, pages 128–137.
- Alshammari, R., Parkes, M., and Adlington, R. (2017). Using whatsapp in efl instruction

- with saudi arabian university students. *Arab World English Journal (AWEJ) Volume*, 8.
- Baadte, C. and Kurenbach, F. (2017). The effects of expectancy-incongruent feedback and self-affirmation on task performance of secondary school students. *Eur. J. Psychol. Educ.*, 32(1):113–131.
- Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). Fitting linear mixed-effects models using lme4. *J Stat Softw*, 67(1):1–48.
- Berge, Z. L. and Collins, M. P. (1995). *Computer mediated communication and the online classroom: distance learning*. Hampton press Cresskill.
- Biber, D. (1992). On the complexity of discourse complexity: A multidimensional analysis. *Discourse Processes*, 15(2):133–163.
- Conrad, S. S. and Dabbagh, N. (2015). Examining the factors that influence how instructors provide feedback in online learning environments. *Int. J. of Online Pedagogy and Course Design*, 5(4):47–66.
- De Smedt, T. (2013). *Modeling Creativity: Case Studies in Python*. University Press Antwerp.
- Dekhinet, R. (2008). Online enhanced corrective feedback for ESL learners in higher education. *Computer Assisted Language Learning*, 21(5):409–425.
- Deshpande, S., Palshikar, G. K., and Athiappan, G. (2010). An unsupervised approach to sentence classification. In *COMAD*, page 88.
- Efron, B. and Tibshirani, R. J. (1994). *An introduction to the bootstrap*. CRC press.
- Graesser, A. C., McNamara, D. S., Cai, Z., Conley, M., Li, H., and Pennebaker, J. (2014). Coh-matrix measures text characteristics at multiple levels of language and discourse. *The Elementary School Journal*, 115(2):210–229.

- Guardado, M. and Shi, L. (2007). ESL students' experiences of online peer feedback. *Computers and Composition*, 24(4):443–461.
- Hannum, W. H., Irvin, M. J., Lei, P.-W., and Farmer, T. W. (2008). Effectiveness of using learner-centered principles on student retention in distance education courses in rural schools. *Distance Education*, 29(3):211–229.
- Harris, L. R., Brown, G. T., and Harnett, J. A. (2015). Analysis of New Zealand primary and secondary student peer-and self-assessment comments: Applying Hattie and Timperley's feedback model. *Assessment in Education: Principles, Policy & Practice*.
- Hattie, J. and Timperley, H. (2007). The power of feedback. *Rev. Educ. Res.*, 77(1):81–112.
- Hepplestone, S., Holden, G., Irwin, B., Parkin, H. J., and Thorpe, L. (2011). Using technology to encourage student engagement with feedback: a literature review. *Research in Learning Technology*, 19(2).
- Ho, M.-c. (2015). The effects of face-to-face and computer-mediated peer review on EFL writers' comments and revisions. *Australasian Journal of Educational Technology*, 31(1).
- Jiang, W. and Ribeiro, A. (2017). Effect of computer-mediated peer written feedback on ESL/EFL writing: A systematic literature review. *Elec. Int. J of Educ., Arts, and Science*, 3(6).
- Kahraman, A. and Yalvac, F. (2015). EFL turkish university students' preferences about teacher feedback and its importance. *Procedia-Social and Behavioral Sciences*, 199:73–80.
- Kimball, L. (2002). Managing distance learning: New challenges for faculty. In *The Digital University Building a Learning Community*, pages 27–40. Springer.

- Lang, C., Heffernan, N., Ostrow, K., and Wang, Y. (2015). The impact of incorporating student confidence items into an intelligent tutor: A randomized controlled trial. *Int. Educ. Data Mining Soc.*
- MacDonald, V. A. (2015). *The application of feedback in secondary school classrooms: Teaching and learning in applied level mathematics*. PhD thesis, University of Toronto (Canada).
- Miller, G. A. (1995). Wordnet: a lexical database for english. *Comm. of the ACM*, 38(11):39–41.
- Miyamoto, Y., Coleman, C., Williams, J., Whitehill, J., Nesterko, S., and Reich, J. (2015). Beyond time-on-task: The relationship between spaced study and certification in MOOCs. *Available at SSRN 2547799*.
- Oinas, S., Vainikainen, M.-P., and Hotulainen, R. (2017). Technology-enhanced feedback for pupils and parents in finnish basic education. *Computers & Education*, 108:59–70.
- Olsen, J. K., Aleven, V., and Rummel, N. (2015). Predicting student performance in a collaborative learning environment. *Int. Educ. Data Mining Soc.*
- Pinto-Llorente, A. M., Sánchez-Gómez, M. C., García-Peñalvo, F. J., and Casillas-Martín, S. (2017). Students’ perceptions and attitudes towards asynchronous technological tools in blended-learning training to improve grammatical competence in english as a second language. *Computers in Human Behavior*, 72:632–643.
- Romero, C. and Ventura, S. (2013). Data mining in education. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 3(1):12–27.
- Saeed, M. A., Ghazali, K., and Aljaberi, M. A. (2018). A review of previous studies on ESL/EFL learners’ interactional feedback exchanges in face-to-face and computer-assisted peer review of writing. *Int. J. of Educ. Tech. in Higher Education*, 15(1):6.

- Senter, R. and Smith, E. A. (1967). Automated readability index. Technical report, CINCINNATI UNIV OH.
- Shang, H.-F. (2017). An exploration of asynchronous and synchronous feedback modes in EFL writing. *J. of Computing in Higher Education*, 29(3):496–513.
- Sinclair, A., McCurdy, K., Lucas, C. G., Lopez, A., and Gašević, D. (2019). Tutorbot corpus: Evidence of human-agent verbal alignment in second language learner dialogues. *Int. Educ. Data Mining Soc.*
- Sinclair, A., Oberlander, J., and Gasevic, D. (2017). Finding the zone of proximal development: student-tutor second language dialogue interactions. *Proceedings of SEMDIAL*, pages 107–115.
- Van der Kleij, F. M., Feskens, R. C., and Eggen, T. J. (2015). Effects of feedback in a computer-based learning environment on students learning outcomes: A meta-analysis. *Rev. Educ. Res.*, 85(4).