# Effective Interfaces for Student-Driven Revision Sessions for Argumentative Writing

Tazin Afrin
University of Pittsburgh
Pittsburgh, USA
tazinafrin@cs.pitt.edu

Omid Kashefi
University of Pittsburgh
Pittsburgh, USA
kashefi@cs.pitt.edu

Christopher Olshefski
University of Pittsburgh
Pittsburgh, USA
cao48@pitt.edu

Diane Litman
University of Pittsburgh
Pittsburgh, USA
litman@cs.pitt.edu

Rebecca Hwa
University of Pittsburgh
Pittsburgh, USA
hwa@cs.pitt.edu

Amanda Godley
University of Pittsburgh
Pittsburgh, USA
agodley@pitt.edu

## ABSTRACT

We present the design and evaluation of a web-based intelligent writing assistant that helps students recognize their revisions of argumentative essays. To understand how our revision assistant can best support students, we have implemented four versions of our system with differences in the unit span (sentence versus sub-sentence) of revision analysis and the level of feedback provided (none, binary, or detailed revision purpose categorization). We first discuss the design decisions behind relevant components of the system, then analyze the efficacy of the different versions through a Wizard of Oz study with university students. Our results show that while a simple interface with no revision feedback is easier to use, an interface that provides a detailed categorization of sentence-level revisions is the most helpful based on user survey data, as well as the most effective based on improvement in writing outcomes.

## CCS CONCEPTS

• **Human-centered computing** → **Interactive systems and tools**; *Graphical user interfaces*; *Web-based interaction*; *Natural language interfaces*; Empirical studies in interaction design; **Empirical studies in HCI**; • **Applied computing** → **Education**; • **Computing methodologies** → *Artificial intelligence*.

## KEYWORDS

Academic writing; revision; argumentative writing; intelligent interface; wizard of oz

## 1 INTRODUCTION

Argumentative writing has long been considered a key component in academic and professional success. Educational research has established that not only does argumentative writing produce positive learning gains among students, but it also contributes to more complex critical thinking skills [7, 15]. However, many students lack the skill of developing an argumentative essay without any writing instruction. Typically, instruction of argumentative writing involves both the composition of multiple drafts of writing and revising those drafts based on formative feedback from others (e.g. teachers, peers). Although most educators and writing instructors agree on the importance of formative feedback, teachers have observed that it can be especially time-consuming, and are thus challenged to consider the balance between efficacy and efficiency [16]. Research on peer feedback suggests that students often do not benefit from peer responses unless peer reviewers have been explicitly instructed how to do it [13].

As a solution, scholars of Natural Language Processing (NLP) have worked toward developing automated writing assistant tools in order to provide instant and constructive feedback to student writers. Many of these tools, however, provide product-focused feedback for one draft at a time (e.g. essay scoring [1], error correction [8], argument mining [4]), as opposed to process-focused feedback, which could provide writers with information not only on the quality of a single draft of writing, but also on the evaluation of their revision patterns from previous to the current draft of an essay. The idea behind ArgRewrite[1], the tool described in this paper, is that improving as a writer involves not only producing increasingly higher quality writing, but it also involves improving on the way one engages in the revision process. The ArgRewrite is designed to help students iteratively revise and update their essays. While previous work shows that feedback on textual revisions encourages students to further revise their essays [30, 33], in this study we want to understand the level of revision categorization (e.g., binary versus detailed) and unit of analysis (sentence or sub-sentential) that is most effective in helping students improve their essay. We hypothesize that a more detailed categorization of a student's revision would be more useful. With that in mind, we design four

---

[1]http://argrewrite.cs.pitt.edu/

web-based interface conditions of the ArgRewrite revision assistant tool – ranging from control with no revision categorization to sentence-level and sub-sentential revision categorization.

This article presents data from a lab-based experiment in which users were provided with one of four different versions of the web-based ArgRewrite tool, each of which differs in unit span of revision analysis and levels of detail in the revision purpose categorization. Condition A is our control interface which provides no feedback at all. Condition B provides binary revision categorization for sentence-level revisions, condition C provides detailed revision categorization for nine different types of sentence-level revisions, and finally condition D used the same revision categorization as C, but provided categorization for sub-sentential revisions. First, we describe the interface components and design decisions for each condition of the ArgRewrite. To understand the usefulness of each condition, we then look at student perception of the system by analyzing the user survey about the interface. Our analysis shows that although our conditions with feedback are not always easy to use compared to the simple control condition, students find the revision categorization helpful to understand their revision effort and weakness. Especially, condition C with detailed sentence-level revision categorization showed to be most useful. Detailed revision categorization also encouraged students to make more revision, qualitatively and quantitatively. We also tested the effectiveness of the system in helping students to further improve their essay score. Again, detailed sentence-level categorization showed to be more useful in helping students boost the essay score. Our research contributions are four fold:

- We developed four conditions of an argumentative revision assistant tool that supports different levels of revision feedback (e.g., binary versus detailed purpose categorization; sentence versus sub-sentential revision unit) and conducted a lab-based study, where students used the tool to revise their essays.
- Using statistical analyses, we compare the usability of the conditions of the tool to understand the revision feedback most helpful from a user perspective.
- Using statistical analyses, we compare the essay score gain to understand what is the best revision feedback to help improve the essay.
- We categorize the revisions students made and perform a comparative analysis to understand the revision behavior by students using different conditions.

## 2 RELATED WORK

Many of the NLP-based writing assistant tools that were developed over the last few years provide feedback on one writing product at a time, or focus on high-level semantic changes. For example, Grammarly [8] provides feedback on grammar mistakes and fluency, ETS-writing-mentor [28] provides feedback to reflect on higher-level essay properties such as coherence, convincingness, etc. Other writing assistant tools such as EliReview [6], Turnitin [22] are designed for peer feedback, plagiarism detection, etc., rather than focusing on writing analysis and feedback. In contrast to those existing tools, we compare two drafts using the ArgRewrite revision assistant tool. While a previous version of ArgRewrite [31] provided feedback based on detailed revision categorization [30, 33] at the sentence-level and was evaluated via a user survey, the current study develops two additional ArgRewrite interfaces (based on binary revision categorization and sub-sentential revision units) and evaluates all interfaces using both user survey and writing improvement analysis.

In terms of revision analysis, work on Wikipedia is the most related to the study of academic writing. Prior works on Wikipedia revision categorization focus on both coarse-level [2] and fine-grained [5, 10, 29] revisions. However, because some fine-grained Wikipedia categories (e.g., vandalism) are specific to wiki scenarios, writing studies instead use fine-grained revision categories more suitable for student argumentative writing [21, 33]. In both cases (Wikipedia or educational), previous studies have focused on investigating the reliability of manually annotating and automatically classifying coarse-level and detailed revision categories, as well as on demonstrating correlations between category frequency and outcome measures. In contrast, our study manipulates whether ArgRewrite provides feedback using coarse-level (surface versus content) or detailed (e.g., claim, evidence, etc.) revision categorizations of textual changes.

Previous studies on writing revision research vary as to whether they use the word-level [2, 5] or the sentence-level as the revision span [31]. Sentences represent a natural boundary of text and automatic revision extraction at the sentence-level has been shown to be reasonably accurate [32]. However, sentence-level revision categories may not always be appropriate. For example, a sentence revision may contain a few fluency changes at the beginning, with substantial information added at the end. In that case, that sentence contains both surface and content revisions. With that in mind, in addition to the sentence-level revisions that were the focus of the original ArgRewrite [31], the current study also explores sub-sentential revisions with detailed revision categorization.

The writer's previous revision effort is often studied in collaborative writing to visualize revisions from multiple authors. For example, DocuViz [25] tracks the number of revisions in google docs and shows the pattern of revising and developing a collaborative document by multiple authors. Unlike collaborative writing, our work focus on multiple revisions by a single author. Another research work that studies visualizing multiple revision patterns by a single student also focuses on the amount of revision through an automated revision graph [17, 18]. Although our ArgRewrite tool does show the number of revisions for each revision category, we do not categorize the revisions based on the frequency. Instead, the revision categories reflect the purpose [33] of that revision. In our tool, the revision are highlighted in both drafts of the essay.

In argument mining, the main goal is to find argument structures and their relations from text. It also focuses on a single text. However, few tools are available for argument mining. One recent work experiments with a text editor to support the student argumentation skills [24]. The tool provides feedback on the argumentation quality of a given text. Students using the tool wrote a more convincing argument than students in the control/baseline condition. A tool called ArguLens helps find issues in issue tracking systems using automatic argument mining [26]. Another recent tool for argument mining is called TARGER [4], which also visualizes argumentative phrases in a text of a single draft. Unlike these argument mining

**Table 1: ArgRewrite interface conditions**

| ArgRewrite Conditions: | A | B | C | D |
|---|---|---|---|---|
| Sentence-level Revision | ✗ | ✓ | ✓ | ✗ |
| Sub-sentence-level Revision | ✗ | ✗ | ✗ | ✓ |
| Binary Revision Categorization | ✗ | ✓ | ✗ | ✗ |
| Detailed Revision Categorization | ✗ | ✗ | ✓ | ✓ |
| Number of participants | 20 | 22 | 22 | 22 |

tools, our ArgRewrite focuses on argumentative revision [33] and compares two drafts of student essays.

Works on formative feedback usually focus on embedded writing instructions for students to further improve the article [11, 19, 27]. While we provide revision analysis and show it with corresponding highlight colors on our web-based tool, this is not a study about providing formative feedback on student essays, or the quality of feedback. Rather, our study focuses on helping students to understand their previous revision effort, or how they addressed the feedback received on the previous draft of an essay. Monitoring one's own progress towards a goal is a cognitively complex task called self-regulation [34, 35]. Previous studies have shown that self-regulation has a positive impact on students' writing development [14, 35]. In our study, self-regulation occurs both during the reflection of previous revision efforts and during the actual revision process. Our ArgRewite tool does not suggest any future revision automatically. Instead, it presents its analysis (but not quality evaluation) of previous revisions so that students can make informed decisions when they further revise the essay.
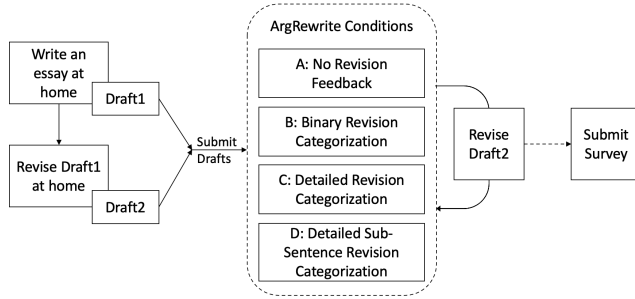
## 3 ARGREWRITE SYSTEM OVERVIEW



**Figure 1: ArgRewrite essay revision process**

Figure 1 shows the essay revision process using the ArgRewrite tool. Experimental participants were recruited through flyers targeting undergraduate and graduate-level students who were either native English speakers or non-native speakers with a certain level of English proficiency (TOEFL score > 100). In our experiment, there are two rounds of essay revision, Draft1 to Draft2, and Draft2 to Draft3. Participants wrote their first draft (Draft1) of an essay at home based on a given prompt[2]. After a few days of finishing Draft1,

[2]The prompt is provided in A.1

each participant received expert feedback[3] on their essay argument quality and overall writing structure. Based on the feedback, they revised their Draft1 and produced Draft2. After finishing Draft2, participants were randomly assigned to use different conditions of the ArgRewrite in a lab environment. They did not receive any feedback on their Draft2. Instead, they are shown the ArgRewrite interface on a computer highlighting their previous revision from Draft1 to Draft2. Participants were asked to use the tool to revise their Draft2 and create a final and generally improved version of the essay, Draft3.

Although our tool supports full automation of revision categorization, we relied on Wizard-of-Oz prototyping [3] for this particular experiment. In Wizard-of-Oz prototyping, a human manually handles the automation, but the student cannot tell the difference from the web-interface they see. We did so to eliminate the confounding factors of NLP automation errors when we compare different conditions. The background server of ArgRewrite uses NLP to automatically segment the essays into sentences and align the two drafts at the sentence-level [31]. Modified, added, or deleted sentences were then extracted as *revisions*. The ArgRewrite server automatically extracts those revisions and classifies them into different revision purpose categories. In our Wizard of Oz experimental setting, a human then fixes the server errors for alignment and classification before the participants start the second round of revision in the lab. In the lab-based experiment, participants first read a short tutorial on using the ArgRewrite tool. Then they were asked to go through their previous revision effort. In conditions B, C, and D, they also submitted confirmation if they agree or disagree with the revision categories for each of the revised sentences the tool is showing them. They did so before and after completing the final revision. Finally, after the participants finished revising the essay, they were asked to answer survey questions about the interface.

Table 1 shows the main differences among the ArgRewrite conditions and the number of participants for each condition. 86 participants were assigned randomly for each condition. Out of 86 participants, 69 were native English speakers, and 17 non-native speakers. The number of non-native speakers in conditions A,B,C,D are 3,4,5,5 respectively. A separate study on participants' native speaking skills showed that non-native speakers made significantly more revisions than native speakers in the first round of revision but not in the second round. Although non-native speakers' scores were lower than native speakers on all drafts and in all conditions, there were no significant differences in non-native vs native speakers revisions or scores across conditions.

## 4 WEB-BASED INTERFACE

Drawing on research on learning analytics [12, 23], ArgRewrite is designed to facilitate personal learning. According to Verbert et al. [23], learning analytics systems provide visualizations and overviews in order to make the users aware of relevant and important information. Each ArgRewrite condition has two parts - the overview interface and the rewrite interface. The overview interface gives a summary of students' revisions between the two

[3]The experts were a professor and a graduate student in the School of Education, and a trained undergraduate student. An example of expert feedback is provided in A.2. The rubric that guided the feedback for Draft1 parallels exactly the scoring rubric (A.3) other than the addition of the pronoun "you."

submitted drafts, while the rewrite interface is where students revise their current draft. Following the previous study [33], in the case of ArgRewrite, the overview interface was designed to bring users' awareness of the purpose of their latest revisions. Then on the rewrite interface, they were asked to go through each revision label to determine whether or not the system identified their revision purposes correctly. Finally, users were allowed to further revise their essay to improve the overall quality.
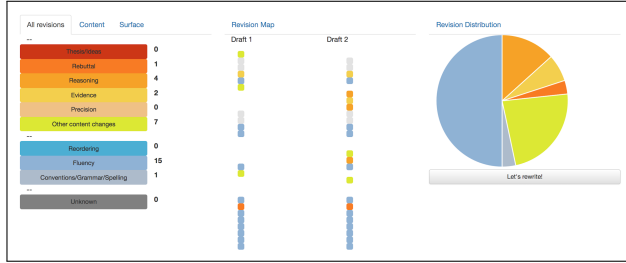
## 4.1 Overview Interface



**Figure 2: Example of the overview interface from ArgRewrite condition C**

The first interface that writers see after logging into ArgRewrite is the Overview interface. Here, writers are presented with overall visualizations of their revision patterns. The three main components of this overview interface are the revision purpose categories, the revision map, and the revision distribution pie chart. Figure 2 shows an example of the overview interface from ArgRewrite condition C. The revision purpose categories are highlighted with their corresponding colors on the left, the revision map is shown in the middle, and the revision distribution pie chart is shown on the right. The components are described below. Once students are ready to revise their essay, they can click on the 'Let's rewrite' button which leads them to the rewrite interface.
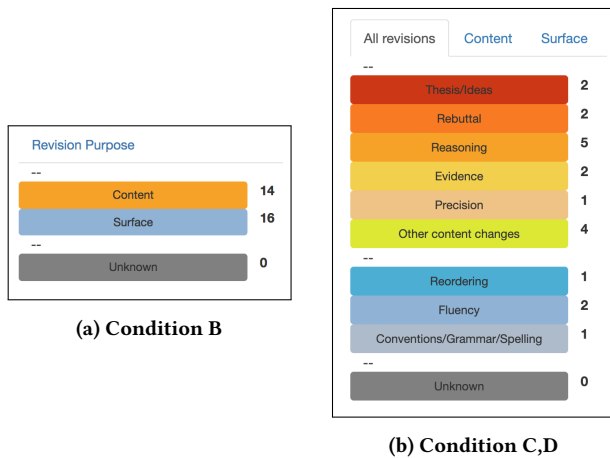


(a) Condition B



(b) Condition C,D

**Figure 3: Revision purpose categories**

*4.1.1 Revision Purpose Categories.* Based on the revision categories presented in [33], our experiment addresses two principal categories of argumentative revisions – surface and content. Surface revisions are the changes that do not alter the meaning of the sentence, e.g., convention or grammar, fluency, and organization changes. Content revisions consist of meaningful textual changes. Following previous works, we use six different categories of content changes – claim, reasoning, evidence, rebuttal, precision, and other general changes[4]. Figure 3 shows the revision purpose categories for different conditions of the ArgRewrite interface. Following previous work [31], surface and content revisions are shown in cold (e.g., blue) and warm (e.g., orange) colors, respectively. Condition B only shows binary revision categories, where the surface and content revisions are shown with blue and orange colors, respectively (shown in Figure 3a). Figure 3b shows the detailed categories and the colors used for conditions C and D. Surface changes in conditions C and D are shown with different levels of blue colors from the cold color scale. Content changes are again shown with warm colors, but take up different colors from the warm color scale. If a revision does not fall into either of those categories, it is labeled as 'unknown' and shown with gray color. The numbers in Figure 3 represent the total added, deleted, and modified revisions for each revision category from Draft1 to Draft2.
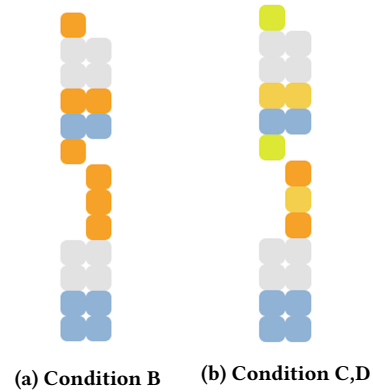


(a) Condition B    (b) Condition C,D

**Figure 4: Revision map shows the drafts as two columns of color-coded tiles, each representing a sentence**

*4.1.2 Revision Map.* Inspired by previous works [20, 31], we design the revision map as two columns of aligned square tiles – the left column represents the previous draft and the right column represents the current draft. Each tile represents a sentence in the draft; the white space between groups of tiles represents the paragraph breaks. Tiles are highlighted with colors of their corresponding revision categories. The shading of the tiles in each row represents whether the student added to, deleted, or modified the original sentence (or made no change). This revision map allows a student to look at all the revisions they made at different locations in the essay at a glance. Students can also easily understand what types of revisions they are making from the highlights. Figure 4 shows the revision map for conditions B, C, and D. In Figure 4a, the first

---

[4]Precision category is added in addition to the content revisions reported in [33].

tile is a deleted sentence because there is no aligned tile/sentence from the current draft. The orange color means it is a content revision. The light gray shade in the next two rows indicates that those sentences are not revised. Tiles in row 4 and 5 indicate modified content and surface revisions respectively. In contrast to the binary categories, Figure 4b shows the same revisions with fine-grained revision categories. It shows that the first sentence is a deleted general content revision, the fourth sentence is modified evidence, and the fifth sentence is a modified fluency revision.
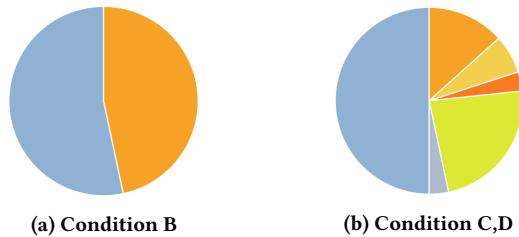


(a) Condition B          (b) Condition C,D

**Figure 5: Revision distribution shown as a pie-chart**

*4.1.3 Revision Pie Chart.* The overview interface also contains a pie chart showing the distribution of the frequency of different revision purpose categories. While the revision purpose categories and the revision map show the number of revisions and the places where revisions are made, the pie chart adds the benefit of easy comparison of the distribution of different types of revisions. Looking at the pie chart, a student can easily understand the influence of the types of revisions they have made between Draft1 to Draft2. Figure 5 shows the revision chart from ArgRewrite conditions B, C, and D. Since we have only two revision types in condition B, Figure 5a shows the distribution of the number of content and surface revisions. This chart (Figure 5a) shows that this student made more surface than content revisions. Figure 5b shows similar information but provides additional details, such as the surface changes were predominately fluency changes, few grammar changes, while the main content changes involved reasoning and other (non-argumentative) content revisions.

## 4.2 Rewrite Interface



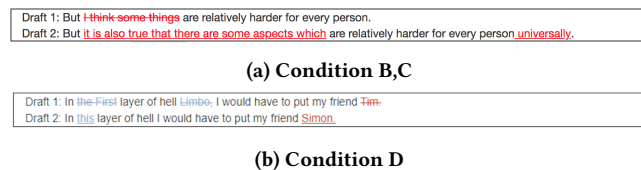(a) Condition B,C



(b) Condition D

**Figure 6: Revision details window for different conditions.**

The rewrite interface contains the revision purpose categories, revision details window, four tabs containing the prompt and three essay drafts, and the revision map similar to the overview interface (except for condition D). Figures 8, 9, and 10 show screenshots of the rewrite interface for different conditions of the ArgRewrite. To encourage students the texts on the drafts tabs are highlighted with

the corresponding revision color. In conditions B and C, the full sentence is highlighted. In condition D, only the revised text within a sentence is highlighted. Students can directly modify the essay on the Draft3 tab, which initially contains Draft2 to start with. When a student clicks on the text to see the details, a small window pops up to show the character-level differences[5] of a selected original and revised sentence. The character differences are highlighted with red in condition B and C. Condition D shows similar differences, but in corresponding revision purpose colors as shown in Figure 6.

The rewrite interface also provides the revision map of sentences to facilitate the navigation through the essay. Students can click on a tile on the revision map on the rewrite interface to look at that particular sentence. However, this is provided for conditions B and C only. Condition D shows a revision map for sub-sentential revisions; it shows two rows of tiles (shown at the top of the Figure 10b) and each tile represents a revised sub-sentential unit within the revised sentences. On the rewrite interface, the small round button beside each tile of the revision map is used to highlight the confirmed revision categories when the students go through their previous revisions and submit their agreement about the revision categories.

## 5 ARGREWRITE CONDITIONS

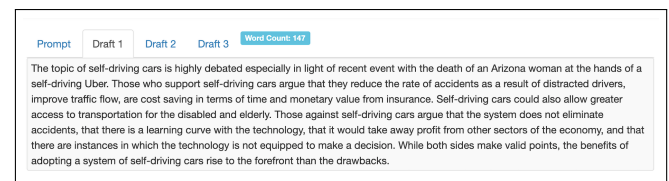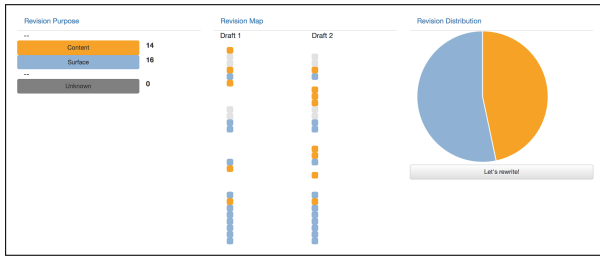## 5.1 Condition A: No Revision Categorization



**Figure 7: A screenshot of the ArgRewrite tool - Condition A**

The ArgRewrite condition A is designed as a baseline containing no revision feedback, to compare with all other ArgRewrite conditions where writers receive different levels of feedback or analysis of their previous revision effort. Since there is no feedback, it does not contain any revision purpose categorization, revision map, or revision pie chart. Therefore, condition A does not have an overview interface. It contains a simplified version of the rewrite interface shown in Figure 7. The rewrite interface contains the plain text of the student essays for each Draft.

## 5.2 Condition B: Binary Revision Categorization

ArgRewrite condition B is designed to provide simple revision feedback to the students. It includes all the components of the overview and the rewrite interface. Revision categorization is shown at the sentence-level. Condition B shows the revisions highlighted using only the top-level (binary) revision purpose categories - surface and content. The surface revisions are highlighted with blue and the content revisions are highlighted with orange to reflect cold versus warm color revisions as described in Section 4.1.1. On the rewrite interface shown in Figure 8b, if a sentence contains any

---

[5]google diff match-patch: https://github.com/google/diff-match-patch

**(a) Overview interface for condition B**



**(b) Rewrite interface for condition B**

**Figure 8: A screenshot of the ArgRewrite tool - Condition B**



**(a) Overview interface for condition C**



**(b) Rewrite interface for condition C**

**Figure 9: A screenshot of the ArgRewrite tool - Condition C**



**(a) Overview interface for condition D**



**(b) Rewrite interface for condition D**

**Figure 10: A screenshot of the ArgRewrite tool- Condition D**

surface revisions, the whole sentence is highlighted with blue. Similarly, sentences with content revisions are highlighted with orange. Similar to condition A, condition B also has four tabs to show the essay prompt and the drafts. Unlike conditions C and D, condition B is simple in terms of categorization of revisions.

## 5.3 Condition C: Detailed Revision Categorization

Condition C shows the detailed revision categorization, highlighted with their corresponding colors shown in Figure 9. It contains all the components of the overview (Figure 9a) and the rewrite interface (Figure 9b). Students get the detailed revision feedback of their essay at sentence-level, according to the revision purpose categories described in Section 4.1.1. In contrast to condition B, students who use condition C to revise their essay can, for example, spot the difference between word-usage versus grammar changes, claim versus evidence changes, etc. It is more informative compared to the control condition and to condition B with its binary revision categorization. Similarly to condition B, the rewrite interface in condition C also shows four tabs and highlights the whole sentence with the identified revision color.
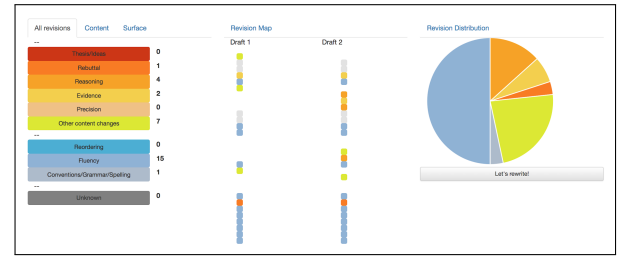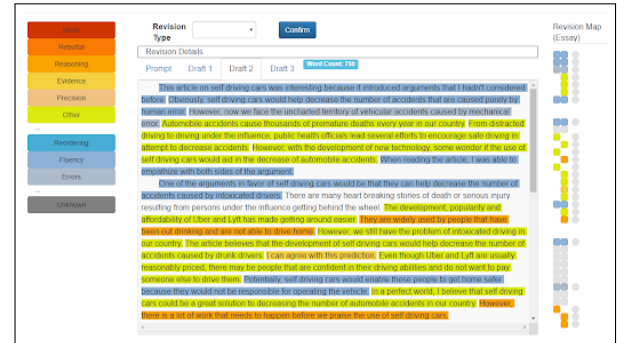
## 5.4 Condition D: Detailed Sub-Sentential Revision Categorization

Condition D is designed to provide more detailed feedback for the revisions students make. Unlike conditions B and C, condition D can focus on multiple different revisions within a single sentence. Each sub-sentential revision is identified and highlighted with the
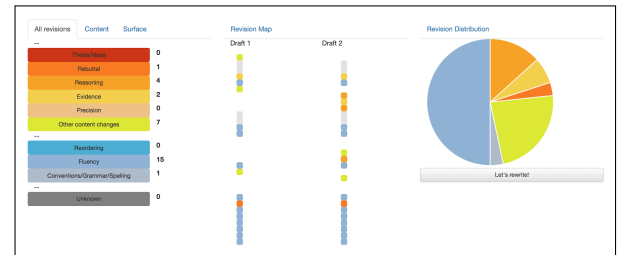
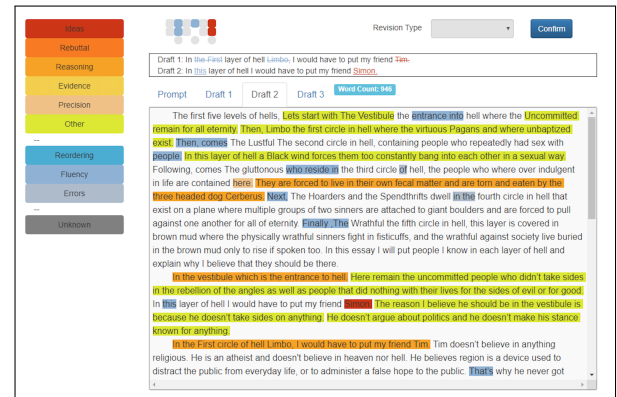corresponding revision category (shown in Figure 10b). This condition contains an overview interface with a sentence-level revision

map, similar to condition C, but the statistics of revision purpose categories are collected and shown from the sub-sentential revision units (Figure 3b and Figure 5b). In the rewrite interface, each sub-sentential revision is highlighted with its corresponding revision purpose color code. By clicking on each revised sentence, a horizontal revision map provides the abstract visualization of how it differs from the original sentence: which sub-sentential units are added, modified, or deleted, and what is the purpose of that revision.

# 6 EVALUATION AND RESULTS

To evaluate our research hypothesis that more detailed feedback is more helpful (i.e., Condition D > C > B > A), we conducted an experiment to answer the following research questions.

**RQ1:** Do students perceive the ArgRewrite to be clear and easy to use?

**RQ2:** Do students find the ArgRewrite helpful for their writing?

**RQ3:** Is ArgRewrite beneficial for student essay improvement?

**RQ4:** Is there any difference in students' revision behavior based on ArgRewrite condition?

Our analyses for RQ1 and RQ2 were based on data from a 16 question survey that participants completed after using ArgRewrite to revise their essays. The survey items addressed [9]'s distinction between "perceived ease of use" and "perceived usefulness" of technology. We included some questions verbatim from [9]'s survey, such as questions 1 and 2, while other items were customized to address unique features of ArgRewrite (shown in Table 2). Eight questions about the perceived ease of use and helpfulness and of the system for supporting essay revision were asked of all participants (questions 1-8). Another set of 8 questions (9-16) focused on usefulness of specific interface components and were asked only of participants in conditions B, C, and D. Each question was answered using a Likert scale ranging from 1 to 5 indicating strongly disagree to strongly agree. To answer RQ3 we examined students' writing improvement, based on expert essay scores that we describe below. Finally, we analyze the revision categories in student essays to answer RQ4. In our analyses, univariate analysis of variance (ANOVA) multiple comparison using Fisher's Least Significant Difference (LSD) test was used to compare differences in survey answers, essay scores, and number of revisions across different conditions. We calculate Cronbach's Alpha coefficient to report internal consistency of the combined survey questions (shown in Table 2). In RQ4, we also use t-test to compare revisions within conditions.

To answer **RQ1**, we combine two survey questions (1-2) that ask about the perceived ease of use of the tool. The questions asked students if they find the system easy to use, and if their interaction with the system is clear and understandable. Mean survey ratings and ANOVA result for those questions are shown in Table 2. For perceived ease of use, the overall difference between conditions is not significant. Looking at pairwise comparison, condition A has a higher mean compared to all other conditions, and Condition D has the lowest mean. Condition A, which is the control condition without any revision feedback, was thus the easiest condition to use. This is not surprising because of the simplicity of the rewrite interface for condition A. However, this mean-value is only significantly higher than condition D, where we provided the most

specific revision feedback. We think this lower mean value reflects the complex information display of the revision categories at the sub-sentence level.

To answer **RQ2**, we first combine the survey questions (3–8) that focus on the perceived usefulness [9] and usage behavior. We then separately examine questions (9–16) regarding usefulness and actual usage of the interface components. Taking the means over questions 3–8 shows that overall, there is a significant difference between conditions although the ANOVA effect size is low. Students perceived condition C with detailed sentence-level revision feedback to be more useful compared to conditions A and B. Particularly, ANOVA results from Table 2 shows that students using condition C thought that the system helps them to better understand their previous revision effort and recognize their weakness, encourages to make more revisions, and more helpful compared to students using conditions A and B. In other words, from this ANOVA result we can say that condition A proved to be less helpful (despite being the easiest to use). Students also perceived detailed sub-sentential revision feedback to be more useful compared to no feedback. For example, when we asked about the quality of revision[6], condition D showed a significantly higher mean-value than condition A. Overall, we can say that detailed feedback is more useful than no feedback or binary feedback which supports our hypothesis. However, we did not see any significant difference between sentence versus sub-sentential revision feedback (C versus D). Therefore we speculate that reducing the granularity of revision feedback might not be very beneficial after all.

We get a mixed signal looking at the questions (9–16) that only target the conditions with feedback (B, C, and D). Overall, ANOVA shows no significant difference between conditions for this group of questions that focus on the actual usage of the interface. However, pairwise comparisons do show some significant differences. For example, students find the revision windows more helpful when they were shown sentence-level revision feedback compared to sub-sentential feedback. However, most of the specific components of the overview and rewrite interface did not show any difference between the conditions (e.g., revision map). On the other hand, a detailed description of revision purpose seemed more inspiring than the binary description (question 10). Detailed sub-sentential feedback was also trustworthy compared to sentence-level binary feedback. Given the Wizard of Oz scenario, the accuracy of the system feedback is objectively similar across conditions.

To answer **RQ3**, we looked at students' essay score. All three drafts written by each participant were scored separately by two researchers, both of whom were experienced high school English and college instructors. The quadratic weighted kappa (QWK) is 0.537. Scoring was guided by a 10-criteria rubric that mirrored the rubric[7] used to give feedback on Draft1 focusing on the argument elements in the essay. Each item was scored on a scale of 1-4: "1-poor," "2-developing," "3-proficient," or "4-excellent." The essay score ranges from 10 to 40. The average of the two researchers' scores was used for data analysis. To determine the improvement of student essay we calculated the normalized essay score gain (NEG) from

---

[6]Students received instruction in the tutorial that content revisions were more related to essay improvement in previous studies [33]. They were encouraged to do more content revisions.

[7]The rubric is provided in A.3

**Table 2: Interface survey questions, mean student response for each condition, and univariate ANOVA result with Fisher's least significant difference (LSD) procedure (* p< .05, ** p< .01, *** p< .001, ~ p< .1, $\alpha$=Cronbach's Alpha)**

| | | A | B | C | D | Pairwise Comparison | F-score | Effect Size |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | ANOVA Result | |
| | Perceived ease of use ($\alpha$=0.87) | 4.40 | 4.16 | 4.09 | 3.75 | D<A* | 2.31~ | 0.08 |
| 1 | I find the system easy to use. | 4.40 | 4.18 | 4.09 | 3.77 | D<A* | 1.69 | 0.06 |
| 2 | My interaction with the system is clear and understandable. | 4.40 | 4.14 | 4.09 | 3.73 | D<A** | 2.48~ | 0.08 |
| | Perceived usefulness and usage behavior ($\alpha$=0.89) | 3.53 | 3.77 | 4.33 | 4.11 | A<C***, A<D**, B<C** | 6.06** | 0.18 |
| 3 | The system allows me to have a better understanding of my previous revision efforts. | 3.70 | 3.95 | 4.45 | 4.27 | A<C**, A<D*, B<C* | 3.40* | 0.11 |
| 4 | The system helps me to recognize the weakness of my essay. | 3.10 | 3.32 | 4.09 | 3.73 | A<C**, A<D*, B<C* | 4.28** | 0.14 |
| 5 | Overall the system is helpful to my writing. | 3.25 | 3.73 | 4.27 | 4.18 | A<C***, A<D***, B<C* | 7.61*** | 0.22 |
| 6 | The system encourages me to make more revisions (quantity) than I usually do. | 3.65 | 3.86 | 4.50 | 4.09 | A<C**, B<C* | 3.28* | 0.11 |
| 7 | The system encourages me to make more meaningful revisions (quality) than I usually do. | 3.45 | 3.86 | 4.23 | 4.23 | A<C*, A<D* | 2.95* | 0.10 |
| 8 | I put a lot of effort into writing and revising this essay. | 4.00 | 3.91 | 4.41 | 4.14 | A<C*, B<C** | 3.91* | 0.13 |
| | Perceived usefulness and actual usage of the interface ($\alpha$=0.68) | – | 4.05 | 4.15 | 4.04 | No difference | 0.41 | 0.01 |
| 9 | I found the overview page to be useful. | – | 4.14 | 4.18 | 4.14 | No difference | 0.03 | 0.00 |
| 10 | The description of the purpose of my revisions inspired me to make more revisions. | – | 3.59 | 4.14 | 4.09 | B<C~ , B<D~ | 2.36 | 0.07 |
| 11 | I found it useful to see my revision purposes highlighted in different colors (i.e. Warm and cold colors) | – | 4.27 | 4.64 | 4.41 | No difference | 1.75 | 0.05 |
| 12 | I found the revision map visualization useful. | – | 4.09 | 3.86 | 3.73 | No difference | 0.72 | 0.02 |
| 13 | I found the small window of revision details to be useful. | – | 4.64 | 4.55 | 3.59 | D<B***, D<C*** | 13.73*** | 0.30 |
| 14 | In general, I found it helpful to know whether my revision was a surface or content level change. | – | 4.05 | 4.09 | 4.23 | No difference | 0.25 | 0.01 |
| 15 | My revision purposes were most often indicated correctly by the system. | – | 4.00 | 3.91 | 3.95 | No difference | 0.08 | 0.00 |
| 16 | I trust the feedback that the system gave me. | – | 3.59 | 3.86 | 4.18 | B < D* | 2.94~ | 0.09 |

**Table 3: Average score and normalized essay score gain (NEG) per condition (* p< .05, ~ p< .1).**

| | | A | B | C | D | Pairwise comparison | F-score | Effect Size |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | ANOVA Result | |
| Score in | Draft1 | 23.38 | 24.15 | 22.80 | 24.23 | | | |
| | Draft2 | 23.90 | 25.64 | 24.55 | 25.25 | - | - | - |
| | Draft3 | 24.90 | 26.50 | 25.84 | 26.11 | | | |
| NEG | 1 to 3 | 0.06 | 0.14 | 0.18 | 0.13 | A<C* , A<B~ | 2.38~ | 0.08 |
| | 2 to 3 | 0.05 | 0.07 | 0.08 | 0.06 | No significant difference | 0.55 | 0.02 |

Draft1 to Draft3 ($NEG13$) and Draft2 to Draft3 ($NEG23$). We did not consider the essay score gain from Draft1 to Draft2 because that step does not involve using our system. Normalized essay score gain is calculated as follows:

$$NEG = \frac{CurrentDraftScore - PreviousDraftScore}{MaxScore - PreviousDraftScore}$$

For both $NEG13$ and $NEG23$, we have the highest mean-value for condition C, where we showed the detailed sentence-level revision feedback (Table 3). We again performed univariate ANOVA with Fisher LSD test to compare the mean of the essay score gains in different interface conditions. The overall ANOVA result did not show any significant difference. ANOVA pairwise comparison result for $NEG13$ showed that students in Condition C performed significantly better than condition A. Condition B was trending better than Condition A ($p = 0.06$). But there was no significant difference between B, C, and D. We also did not see any significant difference for $NEG23$ between any conditions. This result is in line with our previous research question results, in which we observed that students found detailed sentence-level revision feedback to be more helpful compared to no revision feedback at all.

To answer **RQ4**, we looked at the types of revisions (surface vs. content) students made when revising Draft1 to Draft2 (without ArgRewrite) and when revising Draft2 to Draft3 (with ArgRewrite). We expected to see fewer revisions with ArgRewrite since it is the second stage of revising the same essay. Table 4 shows the percentage of surface and content revisions for each condition. Within each condition, we compare the number of surface and content revisions across revision stage using paired t-test. In conditions A and B, we observed significantly more surface revisions and fewer content revisions when revising using ArgRewrite compared to revising without ArgRewrite, but the distribution of types of revisions is not significantly different in condition C and D, when with or without ArgRewrite.

ANOVA result showed no significant difference between conditions for the average number of content or surface revisions. As we have mentioned before, according to previous work, content revisions (e.g. reasoning, evidence) are correlated with essay improvement. Hence, according to Table 4, students in condition A should have higher essay score gains with more content revisions than others. But in Table 3 we have seen that condition A has the lowest essay score gain. With the lowest percentage of content revisions in condition C, students in that condition had higher essay score gains. This result indicates that students who received revision feedback generated revisions that help them improve the essay compared to students who did not receive any feedback. Although students with no feedback generated more content revisions, we speculate those revisions may be irrelevant or unnecessary for supporting the argument.

## 7 DISCUSSION

The findings of this study highlight a tension point that is worth further examination. On the one hand, the analysis of the improvement and revision patterns suggested that Condition C's detailed categorization of revision functions was more effective and helpful than the other conditions. On the other hand, there was an inverse relationship between the granularity of feedback and the usability of the system. In other words, the more detailed the feedback was on students' revision habits, the less students were likely to find it "easy to use" or "clear and understandable" (see questions 1 and 2 on Table 2).

Our findings consistently showed that feedback on detailed revision categorization is better than no feedback. For some evaluation measures, detailed feedback is also better than binary feedback. However, we did not find much difference between sentence versus sub-sentence level revision feedback. So our hypothesis that the more detailed the revision feedback the better is not entirely supported. One potential confound in our study design may have been the different units of analysis employed in Condition D versus the other conditions. By being provided with sub-sentential as opposed to sentential feedback, writers in Condition D spent more time confirming the accuracy of their previous revisions than others. This resulted in them spending more time to look at previous revisions and less time to engage in the actual act of revising when it came to developing their last drafts. This likely contributed to their lower ratings of perceived ease of use, but it also may have influenced the quality of their final drafts. With this in mind, our analyses found little difference between conditions C and D. In the future, we plan to look at the sub-sentence level revisions more closely to understand how to make it more effective for the students. For example, we did not test binary revision categorization at the sub-sentence-level. This is a future condition we would like to explore. Another significant difference we find between sentence-level and sub-sentential interface components is the small window of revision details. Students using sentence-level revision conditions find it more useful than students using sub-sentential revision feedback. We have seen before that the revision details window is different for condition D. It shows the sub-sentence revisions highlighted. So in condition D, students look at the sub-sentential highlights on the essay text and the revision details window, which is redundant. This might be the reason why the revision window was not good enough for condition D but showed to be very useful for conditions B and C.

On one final note regarding our third question related to student improvement, our analyses of improvement from first to third drafts seems to favor detailed sentence-level revision categorization. In our study students revised their Draft1 at home. Hence, the revision from the first to second draft did not involve ArgRewrite. When students used our tool from the second to third draft, they still saw higher essay score gain using sentence-level revision feedback (binary and detailed) than sub-sentential, but those differences were not statistically significant. This might suggest that sub-sentential revision feedback is not helping students improve the essay, even compared to no revision feedback. However, due to the necessary methodological differences mentioned above, we believe we still need to conduct more experiments with sub-sentential revision before reaching any conclusion.

## 8 CONCLUSION

In this paper, we presented a tool that helps students to make further revisions on their argumentative writings. We developed four versions of the interface for the tool and presented a comparative

**Table 4: Number of sentence-level surface and content revisions between first (Draft1 to Draft2) and second (Draft2 to Draft3) revision stage for each condition. (↑: increase in number of revisions compared to the previous revision stage, ↓: decrease in number of revisions compared to the previous revision stage, * p< .05, ~ p< .1)**

| Revision Stage | Revision | A | B | C | D |
|---|---|---|---|---|---|
| 1 to 2 | Surface | 131 (30%) | 136 (37%) | 202 (47%) | 164 (41%) |
| | Content | 322 (70%) | 235 (62%) | 239 (53%) | 230 (59%) |
| 2 to 3 | Surface | 185 (40%) ↑* | 160 (47%) ↑* | 198 (48%) ~ | 183 (45%) ~ |
| | Content | 280 (60%) ↓* | 189 (53%) ↓* | 213 (52%) ~ | 225 (55%) ~ |

study to determine to what extent might the explicit representations of revision purpose categories help students to improve their essay. Our analysis shows that detailed revision categorization at the sentence-level is the most helpful compared to conditions that do not provide detailed feedback. Detailed sub-sentential revision categorization also seemed promising, but more research and development is warranted. In particular, determining the most useful and intuitive level of granularity and detail in writing feedback is an open research question. In the future, we plan to further explore the sub-sentential revision purpose taxonomy to support effective automated writing assistant systems.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Yigal Attali and Jill Burstein. 2006. The Automated Essay Scoring with E-Rater V.2. *Journal of Technology, Learning, and Assessment* 4, 3 (2006).

[2] Amit Bronner and Christof Monz. 2012. User Edits Classification Using Document Revision Histories. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics (EACL '12)*. Association for Computational Linguistics, Avignon, France, 356–366.

[3] Jacob T. Browne. 2019. Wizard of Oz Prototyping for Machine Learning Experiences. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems (Glasgow, Scotland Uk) (CHI EA '19)*. Association for Computing Machinery, New York, NY, USA, 1–6.

[4] Artem Chernodub, Oleksiy Oliynyk, Philipp Heidenreich, Alexander Bondarenko, Matthias Hagen, Chris Biemann, and Alexander Panchenko. 2019. TARGER: Neural Argument Mining at Your Fingertips. In *Proceedings of the 57th Annual Meeting of the Association of Computational Linguistics (ACL'2019)*. Florence, Italy.

[5] Johannes Daxenberger and Iryna Gurevych. 2012. A Corpus-Based Study of Edit Categories in Featured and Non-Featured Wikipedia Articles. In *Proceedings of the 24th International Conference on Computational Linguistics (COLING '12)*. Mumbai, India, 711–726.

[6] The Eli Review. 2014. https://elireview.com. [Online; accessed 01-12-2021].

[7] Jill Fitzgerald and Timothy Shanahan. 2000. Reading and writing relations and their development. *Educational Psychologist* 35, 1 (2000), 39–50.

[8] Grammarly. 2016. http://www.grammarly.com. [Online; accessed 01-12-2021].

[9] Heather Holden and Roy Rada. 2011. Understanding the influence of perceived usability and technology self-efficacy on teachers' technology acceptance. *Journal of Research on Technology in Education* 43, 4 (2011), 343–367.

[10] John Jones. 2008. Patterns of Revision in Online Writing: A Study of Wikipedia's Featured Articles. *Written Communication* 25, 2 (2008), 262–289.

[11] Simon Knight, Antonette Shibani, Sophie Abel, Andrew Gibson, Philippa Ryan, Nicole Sutton, Raechel Wight, Cherie Lucas, Agnes Sandor, Kirsty Kitto, et al.

[12] 2020. AcaWriter: A learning analytics tool for formative feedback on academic writing. *Journal of Writing Research* 12, 1 (2020), 299–344.

[12] Ming Liu, Rafael A. Calvo, and Abelardo Pardo. 2013. Tracer: A Tool to Measure and Visualize Student Engagement in Writing Activities. In *Proceedings of the 2013 IEEE 13th International Conference on Advanced Learning Technologies (ICALT '13)*. IEEE Computer Society, USA, 421–425.

[13] Adam Loretto, Sara DeMartino, and Amanda Godley. 2016. Secondary students' perceptions of peer review of writing. *Research in the Teaching of English* (2016), 134–161.

[14] Charles A MacArthur, Zoi A Philippakos, and Melissa Ianetta. 2015. Self-regulated strategy instruction in college developmental writing. *Journal of Educational Psychology* 107, 3 (2015), 855.

[15] George E. Newell, Jennifer VanDerHeide, and Allison Wynhoff Olsen. 2014. High School English Language Arts Teachers' Argumentative Epistemologies for Teaching Writing. *Research in The Teaching of English* 49 (2014), 95–119.

[16] Trena M. Paulus. 1999. The effect of peer and teacher feedback on student writing. *Journal of Second Language Writing* 8, 3 (1999), 265 – 289.

[17] Antonette Shibani. 2020. Constructing Automated Revision Graphs: A Novel Visualization Technique to Study Student Writing. In *Artificial Intelligence in Education*. Springer International Publishing, Cham, 285–290.

[18] Antonette Shibani, Simon Knight, and Simon Buckingham Shum. 2018. Understanding Revisions in Student Writing Through Revision Graphs. In *International Conference on Artificial Intelligence in Education*. Springer International Publishing, Cham, 332–336.

[19] Valerie J Shute. 2008. Focus on formative feedback. *Review of educational research* 78, 1 (2008), 153–189.

[20] Vilaythong Southavilay, Kalina Yacef, Peter Reimann, and Rafael A. Calvo. 2013. Analysis of Collaborative Writing Processes Using Revision Maps and Probabilistic Topic Models. In *Proceedings of the Third International Conference on Learning Analytics and Knowledge* (Leuven, Belgium) *(LAK '13)*. Association for Computing Machinery, New York, NY, USA, 38–47.

[21] Stephen E. Toulmin. 2003. *The Uses of Argument* (2 ed.). Cambridge University Press.

[22] Turnitin. 2014. http://turnitin.com/. [Online; accessed 01-12-2021]. (2014). http://turnitin.com/

[23] Katrien Verbert, Erik Duval, Joris Klerkx, Sten Govaerts, and José Luis Santos. 2013. Learning Analytics Dashboard Applications. *American Behavioral Scientist* 57, 10 (2013), 1500–1509.

[24] Thiemo Wambsganss, Christina Niklaus, Matthias Cetto, Matthias Söllner, Siegfried Handschuh, and Jan Marco Leimeister. 2020. AL: An Adaptive Learning Support System for Argumentation Skills. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '20)*. Association for Computing Machinery, New York, NY, USA, 1–14.

[25] Dakuo Wang, Judith S. Olson, Jingwen Zhang, Trung Nguyen, and Gary M. Olson. 2015. DocuViz: Visualizing Collaborative Writing. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems* (Seoul, Republic of Korea) *(CHI '15)*. Association for Computing Machinery, New York, NY, USA, 1865–1874.

[26] Wenting Wang, Deeksha Arya, Nicole Novielli, Jinghui Cheng, and Jin L.C. Guo. 2020. ArguLens: Anatomy of Community Opinions On Usability Issues Using Argumentation Models. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '20)*. Association for Computing Machinery, New York, NY, USA, 1–14.

[27] Ursula Wingate. 2010. The impact of formative feedback on the development of academic writing. *Assessment & Evaluation in Higher Education* 35, 5 (2010), 519–533.

[28] The Writing Mentor. 2016. ETS Writing Mentor, https://mentormywriting.org/, [Online; accessed 01-12-2021].

[29] Diyi Yang, Aaron Halfaker, Robert E. Kraut, and Eduard H. Hovy. 2017. Identifying Semantic Edit Intentions from Revisions in Wikipedia. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP'17)*.

Association for Computational Linguistics, Copenhagen, Denmark, 9–11.

[30] Fan Zhang, Homa Hashemi, Rebecca Hwa, and Diane Litman. 2017. A Corpus of Annotated Revisions for Studying Argumentative Writing. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (Vancouver, Canada). Association for Computational Linguistics, 1568–1578.

[31] Fan Zhang, Rebecca Hwa, Diane Litman, and Homa B. Hashemi. 2016. ArgRewrite: A Web-based Revision Assistant for Argumentative Writings. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*. Association for Computational Linguistics, San Diego, California, 37–41.

[32] Fan Zhang and Diane Litman. 2014. Sentence-level Rewriting Detection. In *Proceedings of the Ninth Workshop on Innovative Use of NLP for Building Educational Applications*. Association for Computational Linguistics, Baltimore, Maryland, 149–154.

[33] Fan Zhang and Diane Litman. 2015. Annotation and Classification of Argumentative Writing Revisions. In *Proceedings of the 10th Workshop on Innovative Use of NLP for Building Educational Applications*. Association for Computational Linguistics, Denver, Colorado, 133–143.

[34] Barry Zimmerman and Anastasia Kitsantas. 2002. Acquiring writing revision and self-regulatory skill through observation and emulation. *Journal of Educational Psychology* 94 (12 2002), 660–668.

[35] Barry J. Zimmerman and Albert Bandura. 1994. Impact of Self-Regulatory Influences on Writing Course Attainment. *American Educational Research Journal* 31, 4 (1994), 845–862.

# A  DATA COLLECTION MATERIALS

## A.1  Prompt

In this argumentative writing task, imagine that you are writing an op-ed piece for the Pittsburgh City Paper about self-driving cars. The editor of the paper has asked potential writers, like you, to gather information about the use of self-driving cars, and argue whether they are beneficial or not beneficial to society. In your writing, first, briefly explain both the advantages and disadvantages of self-driving cars. Then, you will choose a side, and construct an argument in support of self-driving cars as beneficial to society, or against self-driving cars as not beneficial to society.

A high quality op-ed piece maintains a clear position on the issue and uses supporting ideas, strong evidence from the reading, explanations of your ideas and evidence, and a counter-argument. Furthermore, a high quality op-ed piece is clearly organized, uses precise word choices, and is grammatically correct.

## A.2  Example of Expert Feedback on Draft1

Thank you for your participation in the study. Your draft has been read, and feedback from an expert writing instructor is written below. We advise that you use this feedback when you revise.

The strengths of your essay include:

- All claims have relevant supporting evidence, though that evidence may be brief or general.
- You respond to one, but not all parts of the prompt. However, your entire essay is focused on the prompt.

Areas to improve in your essay include:

- You provided a statement that somewhat show your stance for or against self-driving cars, but it is unclear, or is just a restatement of the prompt.
- Your essay's sequence of ideas is inconsistent, with some clear and some unclear progression.
- Your essay does not include a rebuttal.

## A.3  Scoring Rubric

Table 5 shows the scoring rubric used to provide feedback.

## Table 5: Argumentative Essay Rubric

| | 1-Poor | 2-Developing | 3-Proficient | 4-Excellent |
|---|---|---|---|---|
| **Response to prompt** | The essay is off topic, and does not consider or respond to the prompt in any way. | The essay addresses the topic, but the entire essay is not focused on the prompt. The author may get off topic at points. | The author responds to one, but not all parts of the prompt, but the entire essay is focused on the prompt. | The author responds to all parts of the prompt and the entire essay is focused on the prompt. |
| **Thesis** | The author did not include a statement that clearly showed the author's stance for or against self-driving cars. | The author provided a statement that somewhat showed the author's stance for or against self-driving cars, though it may be unclear or only a restatement of the essay prompt. | The author provided a brief statement that reflects a thesis, and is indicative of the stance the author is taking toward self-driving cars. | The author provided a clear, nuanced and original statement that acted as a specific stance for or against self-driving cars. |
| **Claims** | The author's claims are difficult to understand or locate. | The author's claims are present, but are unclear, not fully connected to the thesis or the reading, or the author makes only one claim multiple times. | The author makes multiple, distinct, and clear claims that align with either their thesis or the given reading, but not both. | The author makes multiple, distinct claims that are clear, and align with both their thesis statement and the given reading. They fully support the author's argument. |
| **Evidence for Claims** | The author does not provide any evidence to support thesis/claims. | Less than half of claims are supported with relevant or credible evidence or the connections between the evidence and the thesis/claims is not clear. | All claims have relevant supporting evidence, though that evidence may be brief or general. The source of the evidence is credible and acknowledged/cited where appropriate. | The author provides specific and convincing evidence for each claim, and most evidence is given through detailed personal examples, relevant direct quotations, or detailed examples from the provided reading. The source of the evidence is credible and acknowledged/cited where appropriate. |
| **Reasoning** | The author provides no reasoning for any of their claims. | Less than half of claims are supported with reasoning or the reasoning is so brief, it essentially repeats the claim. Some reasoning may not appear logical or clear. | All claims are supported with reasoning that connect the evidence to the claim, though some may not be fully explained or difficult to follow. | All claims are supported with clear reasoning that shows thoughtful, elaborated analysis. |

| | 1-Poor | 2-Developing | 3-Proficient | 4-Excellent |
|---|---|---|---|---|
| **Reordering/ Organiza- tion** | The sequence of ideas/claims is difficult to follow and the essay does not have an introduction, conclusion, and body paragraphs that are organized clearly around distinct claims. | The essay's sequence of ideas is inconsistent, with some clear and some unclear progression of ideas OR the essay is missing a distinct introduction OR conclusion. | The essay has a clear introduction, body, and conclusion and a logical sequence of ideas, but each claim is not located in its own separate paragraph. | The essay has an introduction, body and conclusion and a logical sequence of ideas. Each paragraph makes a distinct claim. |
| **Rebuttal** | The essay does not include a rebuttal. | The essay includes a rebuttal in the sense that it acknowledges another point of view, but does not explore possible reasons why this other viewpoint exists. | The essay includes a rebuttal in the form of an acknowledgement of a different point of view and reasons for that view, but does not explain why those reasons are incorrect or unconvincing. | The essay explains a different point of view and elaborates why it is not convincing or correct. |
| **Precision** | Throughout the essay, word choices are overly informal and general (e.g., "I don't like self-driving cars because they have problems."). | Word choices are mostly overly general and informal, though at times they are specific. | Word choices are mostly specific though there may be a few word choices that make the meaning of the sentence vague. | Throughout the essay, word choices are specific and convey precise meanings (e.g., "Self-driving cars are dangerous because the technology is still not advanced enough to address the ethical decisions drivers must make.") |
| **Fluency** | A majority of sentences are difficult to understand because of incorrect/ inappropriate word choices and sentence structure. | A noticeable number of sentences are difficult to understand because of incorrect/ inappropriate word choices and sentence structure, although the author's overall point is understandable. | Most sentences are clear because of correct and appropriate word choices and sentence structure. | All sentences are clear because of correct and appropriate word choices and sentence structure. |
| **Conventions/ Grammar/ Spelling** | The author makes many grammatical or spelling errors throughout their piece that interfere with the meaning. | The author makes many grammatical or spelling errors throughout their piece, though the errors rarely interfere with meaning. | The author makes few grammatical or spelling errors throughout their piece, and the errors do not interfere with meaning. | The author makes few or no grammatical or spelling errors throughout their piece, and the meaning is clear. |