

# Direct Automated Feedback Delivery for Student Submissions based on LLMs

MAXIMILIAN SÖLCH, Technical University of Munich, Germany

FELIX T.J. DIETRICH, Technical University of Munich, Germany

STEPHAN KRUSCHE, Technical University of Munich, Germany

Timely and individualized feedback is crucial for students' learning progress and motivation. However, providing such feedback is a major challenge in education, especially as the number of students has steadily increased in recent years. This growth has made it difficult for teachers to provide individualized feedback to each student, resulting in a time-consuming, repetitive and often manual task that contributes to a high workload.

This paper presents Direct Automated Feedback (DAFeD), a large language model (LLM)-based approach for automated formative feedback on student submissions in various exercise domains. The defined feedback process enables interactive learning by allowing students to submit their solutions multiple times and request feedback repeatedly before the submission deadline. By incorporating task-specific information into the prompt, DAFeD provides iterative customized feedback, facilitating continuous student improvement.

To empirically evaluate the feedback process, we implemented it in an open-source reference implementation that is integrated into the LMS learning platform. We conducted a controlled study in which students used the feedback process in a programming task in a supervised environment and completed a survey. The results show that students perceive the feedback process as relevant and beneficial to their learning. Students indicated that they feel more comfortable and willing to request automated feedback than they would with human tutors because it is more convenient and immediate. This shows that DAFeD has the potential to significantly improve the feedback process in educational institutions and increase students' learning efficiency and performance.

CCS Concepts: • **Social and professional topics** → **Student assessment**; • **Applied computing** → **Education**.

Additional Key Words and Phrases: Software Engineering, Education, Formative Feedback, Interactive Learning

## ACM Reference Format:

Maximilian Sölch, Felix T.J. Dietrich, and Stephan Krusche. 2024. Direct Automated Feedback Delivery for Student Submissions based on LLMs. In . ACM, New York, NY, USA, 11 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

In the current educational landscape, providing timely and effective feedback to students remains a significant challenge. Traditionally, students must wait for course tutors or professors to review their submissions and provide feedback. This process can be time-consuming, often requiring students to arrange meetings and wait for available time slots, which are not always convenient or immediate. Similar it is timeconsuming and tedious for professors and tutors to provide asynchronous feedback via email or other communication channels. The inherent delays and scheduling difficulties make this approach not scalable, especially in courses with a large number of students.

These limitations hinder students' learning progress and motivation. The waiting period for feedback interrupts the learning flow, causing students to lose momentum and potentially disengage from the subject matter. Additionally, the limited availability of tutors and professors means that not all students receive the individualized attention they need

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

to improve their understanding and skills. This situation underscores the necessity for a more efficient and scalable feedback system that can provide continuous support to students without the constraints of traditional methods.

In this paper, we present an approach for generating automated feedback on student submissions using large language models (LLMs) to address these challenges. The approach is independent of the exercise type and can be applied to various domains, such as programming, text, or modeling exercises. We implemented the approach in an open-source reference implementation called Athena, connected to the learning platform Artemis through which students submit their solutions and receive feedback. To validate the effectiveness and efficiency of the approach, we tested it in a controlled environment. We collected quantitative and qualitative data to evaluate students' perceptions of the approach and the overall performance of the reference implementation. The results show

The subsequent sections of this paper are organized to provide a comprehensive understanding of the research. Section 2 provides an overview of related work. Section 3 details the concept and methodology of Direct Automated Feedback Delivery (DAFeeD). Section 4 describes the reference implementation of DAFeeD, called Athena, including a general overview, details on the used prompts, and the system architecture. Section 5 presents the evaluation results. Finally, Section 6 concludes with a summary of findings and discusses future research directions to enhance automated feedback systems.

## 2 RELATED WORK

The necessity and challenges of providing effective feedback in educational settings have been widely recognized. This section overviews prior work in three interrelated areas: automated feedback systems, the application of large language models in educational feedback, and feedback delivery mechanisms within learning platforms.

### 2.1 Automated Feedback Systems

todo [2]

### 2.2 Large Language Models in Educational Feedback

### 2.3 Feedback Delivery in Learning Platforms

### 3 APPROACH: DIRECT AUTOMATED FEEDBACK DELIVERY (DAFEED)

DAFeeD employs large language models to deliver automated feedback on student submissions, designed to complement traditional teaching methods and provide additional support. Figure 1 illustrates the continuous feedback workflow that DAFeeD facilitates, enabling students to receive feedback at any time, thereby eliminating the need to wait for responses from human tutors or course professors.

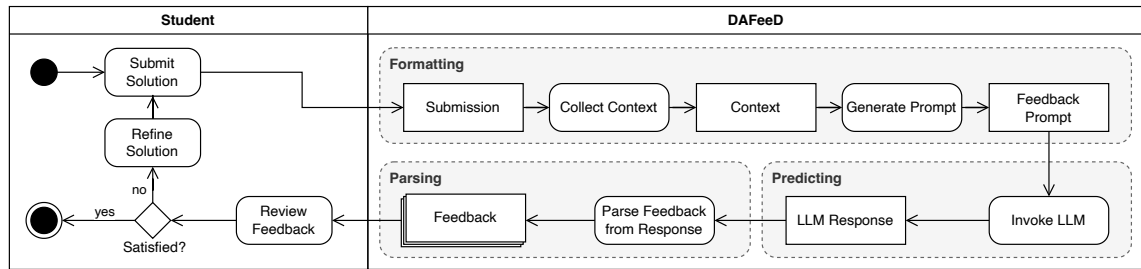


Fig. 1. Workflow of direct automated feedback delivery (DAFeeD) for students' submissions (UML Activity Diagram)

The feedback process is designed to be exercise-independent, meaning that it can be applied to various exercise types, such as programming, text, or modeling exercises. DAFeeD can provide formative feedback to the students, including feedback on issues or improvements, as well as positive feedback when the student completes the task correctly. Once the student submits their solution, DAFeeD initiates a three-stage process to generate natural language feedback.

The first stage, called *Formatting*, takes the student's submission and extracts the submission content, the problem statement, including learning objectives, and any possible grading instructions the instructor defines. All of this gathered information represents the context. During the prompt generation step, a predefined prompt template is filled with the prompt input data, resulting in the feedback prompt. Depending on the exercise, adaptations need to be made to the prompt template to ensure that the feedback output of the LLM is tailored to the specific exercise type. For programming exercises, the generated feedback needs to have metadata information about the file and line number of the code snippet to which the feedback refers. For text exercises, the feedback needs to have metadata about the sentence or word range the feedback refers to.

In the second stage, called *Predicting*, DAFeeD sends the feedback prompt to a large language model (LLM) and invokes it with the prompt. As a result, the LLM generates a response to that prompt including detailed feedback items for the student.

The final stage, *Parsing*, takes the LLM response, which comes in the JSON format, and parses feedback items from it. In addition to the feedback text, the feedback object also contains reference information indicating the part of the submission it pertains to. For programming exercises, this includes the file name and line number of the relevant code snippet to which the feedback refers. For text exercises, the reference information includes only the sentence or word range the feedback refers to.

All of the feedback is then returned to the student for review. If the student is satisfied with the feedback, the process concludes. Otherwise, the student can refine and resubmit their solution, initiating the DAFeeD process anew.

This iterative process is designed to motivate students to continuously learn and experiment with their solutions, resulting in improved performance.

## 4 REFERENCE IMPLEMENTATION: ATHENA

We incorporated DAFeeD into a reference implementation named Athena, which is seamlessly integrated with the learning platform Artemis. Through Artemis, students can submit their solution and review the feedback.

When submitting their solutions on Artemis, students have the option to request direct automated feedback by clicking a newly added button. This feedback request is then sent to Athena, provided the student has not reached their feedback request limit for the exercise. Course instructors can customize the number of allowed feedback requests per exercise according to their preference. A status visualization informs students about their feedback request state. Once Athena generates the feedback and sends it back to Artemis, the student can review it in a modal window on Artemis, as depicted in Figure 2.

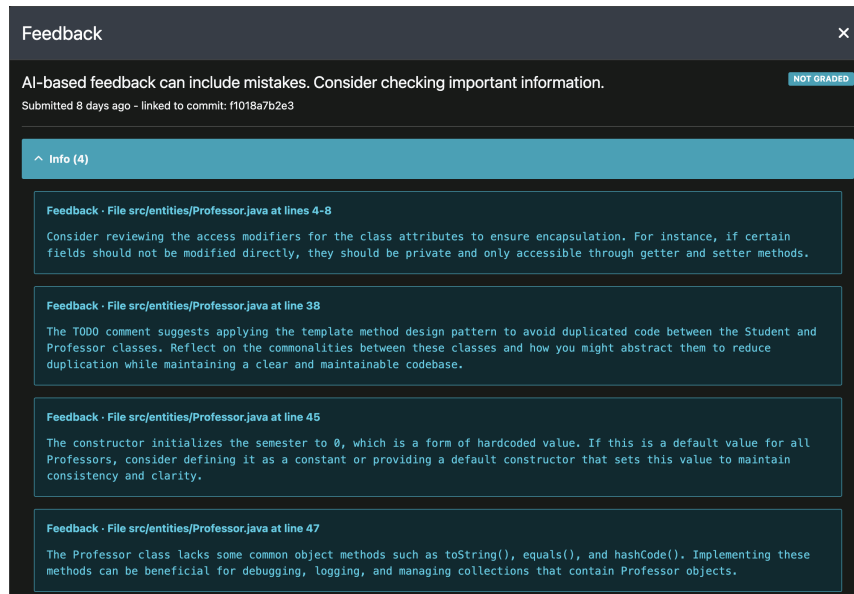


Fig. 2. Visualization of the feedback how students see it in Artemis.

### 4.1 Feedback Generation

The prompt design is crucial for guiding the large language model (LLM) in generating effective and contextually relevant feedback. In Figure 3, we provide an example of a prompt used for generating feedback for programming exercises. This prompt incorporates specific instructions to ensure that the feedback is tailored to the student's submission.

The feedback generation process begins by identifying the differences between the student's submission repository and the provided template repository, the starting point. These differences are identified using a git diff, which highlights lines removed and added by the student. If the problem statement is too lengthy or complex, a separate LLM invocation is used to split the problem statement into relevant parts for each file. This ensures that the feedback is targeted and relevant to the specific context of the file being reviewed. Additionally, a summary of the student's solution across all files is generated using another LLM invocation. This summary provides a comprehensive overview of the submission, which is included in the prompt to offer context for the feedback.

In the provided prompt, several key components guide the AI in creating useful feedback. The *Problem Statement* section contextualizes the student's task and helps the AI understand the exercise's objectives. The *Task Instructions* direct the AI to provide non-graded, constructive feedback focusing on educational aspects without offering direct solutions. *Style Guidelines* ensure the feedback is constructive, specific, balanced, clear, concise, actionable, educational, and contextual. The *File Path and Content* provide the specific file under review along with its content, aiding the AI in pinpointing specific lines of code for feedback. Additionally, *Summary and Diffs* between the template and submission offer additional context, helping the AI understand the student's changes and their overall approach.

The structure and content of this prompt are designed to emulate a human tutor's approach, ensuring that the feedback is both relevant and supportive of the student's learning process. By providing such detailed instructions and contextual information, the LLM can generate feedback that is both meaningful and actionable for students.

```
You are an AI tutor for programming assessment at a prestigious university .

# Problem statement
{problem_statement}

# Task
Provide constructive , non-graded feedback on a student 's programming submission as a
human tutor would. The tutor is not familiar with the solution , so the feedback should
focus solely on aspects from which the student can learn. This feedback must highlight
incorrectly applied principles or inconsistencies without offering specific solutions or
error corrections. Allow some flexibility for students to deviate from the problem
statement , provided they complete all tasks. Ensure the feedback is balanced and
comprehensive .

# Style
1. Constructive , 2. Specific , 3. Balanced , 4. Clear and Concise , 5. Actionable , 6.
Educational , 7. Contextual

Feedback that contradicts the problem statement is strictly prohibited . Avoid mentioning
aspects not explicitly covered in the template to submission diff , such as the exercise
package name , as these are beyond the student 's control .

In git diff , lines marked with '-' were removed and with '+' were added by the student .

The student will be reading your response , so use "you" instead of "them" .

Path : {submission_file_path}

File (with line numbers <number>: <line>) :
{submission_file_content}

Summary of other files in the solution :
{summary}

The template to submission diff (only as reference) :
{template_to_submission_diff}
```

Fig. 3. Prompt for Generating Feedback for Programming Exercises. The highlighted sections are placeholders for the respective elements.

## 4.2 Architecture

Athena is deployed in production alongside the learning platform Artemis, which serves up to more than 2000 students per course. Consequently, the reference implementation must satisfy additional non-functional requirements such as performance, scalability, maintainability, and usability. To meet these requirements and to support feedback generation for multiple exercise types while allowing for future extensibility, we adopted a modular architecture, as illustrated in Figure 4.

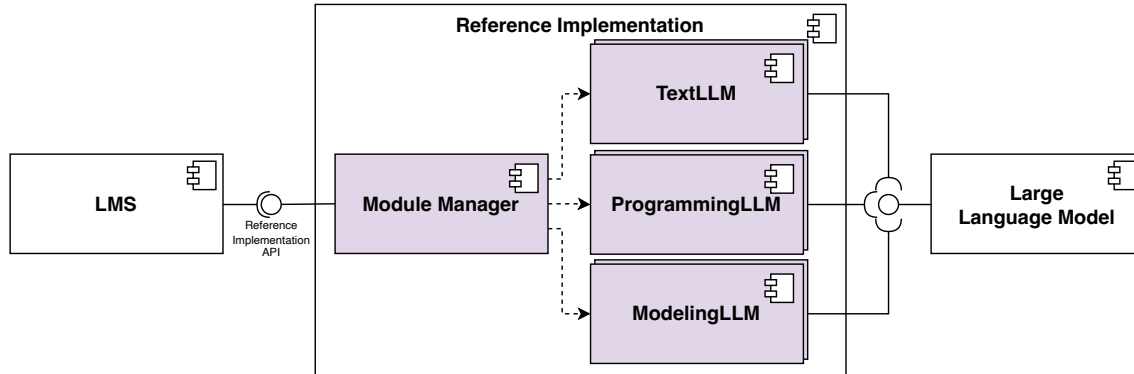


Fig. 4. Top Level Architecture of the reference implementation Athena (UML Deployment Diagram)

The *Module Manager* handles all incoming requests, verifies authorization, and forwards them to the appropriate modules. The *ProgrammingLLM* module manages programming exercises and executes the three-stage DAFeeD process, which includes formatting, predicting, and parsing. Similarly, the *TextLLM* module is optimized for text exercises and follows the same process.

Athena's system design is independent of any specific learning management system (LMS) as it provides a REST API, documented using the OpenAPI standard<sup>1</sup>. This independence allows Athena to be integrated with various LMS platforms, such as Moodle<sup>2</sup>.

Athena currently connects to OpenAI models hosted in a private Azure cloud to ensure that student data is not used for training models, maintaining privacy. Additionally, the system can be configured to use open-source models like Llama<sup>3</sup> or Mistral<sup>4</sup>, either self-hosted or cloud-based.

To meet performance and scalability requirements, Athena and its modules are deployed within a Kubernetes cluster<sup>5</sup>. Kubernetes, in conjunction with Athena's modular architecture, allows the system to scale each module independently. For example, additional instances of the programming module can be instantiated when a new programming exercise is released. Furthermore, Kubernetes provides out-of-the-box load balancing and self-healing capabilities, ensuring that if a module crashes, it is automatically restarted.

<sup>1</sup><https://www.openapis.org>

<sup>2</sup><https://moodle.org>

<sup>3</sup><https://llama.meta.com>

<sup>4</sup><https://mistral.ai>

<sup>5</sup><https://kubernetes.io>

## 5 EVALUATION

In this section, we outline the methodology employed to validate the effectiveness of the proposed DAFeeD approach including the reference implementation Athena. The conducted evaluation represents the treatment validation stage of the design science methodology proposed by Wieringa [4]. In this stage, the proposed solution — DAFeeD — is evaluated in a controlled environment, and the collected data is utilized for the refinement and improvement of the solution.

We begin by describing the research questions, followed by the study design and the results. Subsequently, we outline the limitations of the evaluation and discuss the implications of the findings.

### 5.1 Research Questions

With this study, we want to answer the following research questions about direct automated feedback delivery:

**RQ1** How does the availability of direct automated feedback affect student engagement and motivation?

**RQ2** Do students feel more comfortable requesting automatic feedback than asking a human tutor or the course professor?

**RQ3** How do students perceive the effectiveness of direct automated feedback?

**RQ4** How do students perceive the usability and helpfulness of DAFeeD?

RQ1 examines the influence of direct automated feedback on students' overall engagement with the course material and their motivation to complete exercises and improve their skills. We aim to determine if the immediacy and convenience of automated feedback enhance students' commitment to their coursework. RQ2 explores students' comfort levels with seeking feedback from an automated system compared to traditional sources such as human tutors or course professors. We want to assess whether students prefer the anonymity and immediacy of automated feedback over potentially intimidating interactions with instructors. RQ3 seeks to understand students' views on the value and impact of the feedback provided by the DAFeeD system. We are interested in whether students find the feedback to be relevant and beneficial to their learning process. RQ4 focuses on students' perceptions of the usability and helpfulness of the DAFeeD system. We aim to evaluate how intuitive and user-friendly students find the system, as well as how effective they consider it in assisting their learning and improving their performance.

### 5.2 Study Design

To gain comprehensive insights into students' perceptions of the newly introduced DAFeeD concept, we employed a survey-based approach. Initially, study participants tested the new feedback feature on a sample exercise using the Artemis platform in a controlled environment at the university. The study conductors were present throughout to answer questions and provide support as needed. Following this hands-on experience, participants were asked to complete a survey hosted on the community version of the open-source survey tool LimeSurvey<sup>6</sup>. This survey aimed to gather their opinions on direct automated feedback and collect feedback on their overall experience with the feature.

We invited students from current courses at the university to participate in the study via direct messages and conducted the study with a total of 20 participants. The study uses a mixed methods approach, combining quantitative and qualitative data collection methods. The participants were a mix of undergraduate and graduate students from various disciplines, including computer science, information systems, and games engineering.

<sup>6</sup><https://www.limesurvey.org>



All survey questions, except for the introductory demographic queries and the final, voluntary free-text responses, employ a 5-point Likert scale [1] ranging from "strongly agree" to "strongly disagree" and are mandatory. The following list shows the mapping from the asked questions and statements to the research questions:

#### RQ1 Engagement and Motivation

- Q1 The direct automated feedback keeps me more engaged in the learning process.
- Q2 The direct automatic feedback motivates me to repeatedly improve my code.
- Q3 The direct automated feedback makes me feel more motivated to complete my programming assignments.
- Q4 The direct automated feedback encourages me to experiment more with my coding solutions.

#### RQ2 Comfort with Feedback Source

- Q5 I feel more comfortable requesting direct automated feedback than feedback from a human tutor.
- Q6 I am likely to request feedback more frequently when using direct automated feedback than feedback from my course professor.
- Q7 I find receiving direct automated feedback less intimidating than receiving feedback from a human tutor.
- Q8 I feel that requesting direct automated feedback is more convenient than arranging a meeting with a human tutor.

#### RQ3 Perceived Effectiveness

- Q9 The direct automated feedback helps me understand my mistakes.
- Q10 The direct automated feedback is more effective than one-time feedback.
- Q11 The direct automated feedback has significantly improved the quality of my programming assignment.
- Q12 The direct automated feedback is a helpful addition to the automatic test case results.
- Q13 I feel that having access to direct automated feedback continuously helps me more than arranging a meeting with a human tutor.

#### RQ4 Usability and Helpfulness

- Q14 It is easy to receive direct automated feedback on my programming assignments.
- Q15 I would rather use the direct automated feedback integrated into Artemis than use an external AI tool for getting feedback.
- Q16 I find the direct automated feedback helpful in improving my programming skills.
- Q17 I am satisfied with the overall performance of the direct automated feedback.

### 5.3 Results

In the following paragraphs, we present the results. The answers to each of the likert scale questions are visualized in Figure 5.

90% of students indicated that the direct automated feedback from Artemis keeps them more engaged in the learning process, with 10% neutral (Q1). For Q2, 85% of students stated that the direct automatic feedback from Artemis motivates them to repeatedly improve their code, with 10% neutral and 5% disagreeing. For Q3, 85% of students mentioned that the direct automated feedback from Artemis makes them feel more motivated to complete their programming assignments, with 10% neutral and 5% disagreeing. The encouragement to experiment more with coding solutions (Q4) received a positive response, with 70% agreeing, 20% neutral, and 10% disagreeing.

Comfort levels in requesting feedback showed that 90% of students feel more comfortable requesting direct automated feedback from Artemis than feedback from a human tutor, with 10% neutral (Q5). For Q6, 100% of students noted that

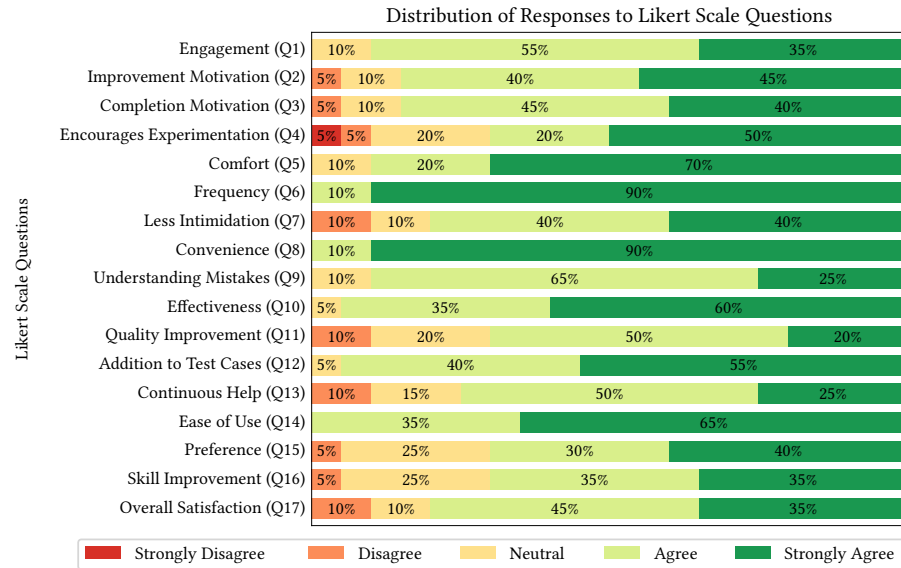


Fig. 5. Distribution of Likert scale responses.

they are likely to request feedback more frequently when using direct automated feedback from Artemis than feedback from their course professor. Receiving feedback from Artemis was found to be less intimidating than from a human tutor by 80% of students, with 10% neutral and 10% disagreeing (Q7). Convenience in requesting feedback showed that 100% of students feel that requesting direct automated feedback from Artemis is more convenient than arranging a meeting with a human tutor (Q8).

In terms of understanding mistakes, 90% of students believed that the direct automated feedback provided by Artemis helps them understand their mistakes, with 10% neutral (Q9). The effectiveness of the feedback was highlighted by 95% of students who found that the direct automated feedback from Artemis is more effective than one-time feedback, with 5% neutral (Q10). Regarding the quality of assignments, 70% of students observed that the direct automated feedback has significantly improved the quality of their programming assignments, with 20% neutral and 10% disagreeing (Q11). For Q12, 95% of students felt that the direct automated feedback is a helpful addition to the automatic test case results, with 5% neutral. Continuous access to feedback was found to be more beneficial than arranging meetings with a tutor by 75% of students, with 15% neutral and 10% disagreeing (Q13).

Ease of receiving feedback was highly rated, with 100% of students confirming that it is easy to receive direct automated feedback from Artemis on their programming assignments (Q14). Furthermore, 70% of students preferred using the direct automated feedback integrated into Artemis over using an external AI tool for getting feedback, with 25% neutral and 5% disagreeing (Q15). In terms of skill improvement, 70% of students agreed that they find the direct automated feedback from Artemis helpful in improving their programming skills, with 25% neutral and 5% disagreeing (Q16). Lastly, 80% of students were satisfied with the overall performance of the direct automated feedback system, with 10% neutral and 10% disagreeing (Q17).

The responses to the voluntary free text questions highlight several themes. Many students appreciated the immediate availability of feedback, which allowed for prompt corrections without waiting for manual review. However, some

respondents suggested improvements such as better categorization of feedback, more detailed explanations of errors, and prioritization of critical issues. The feedback was generally found to be relevant and useful in addressing obvious mistakes and improving code quality. Students expressed a preference for feedback that clearly identified mistakes and provided specific guidance on how to correct them, along with suggestions for improvement. Some challenges included understanding certain automated feedback messages and occasional false positives or negatives in error detection.

#### 5.4 Findings

#### 5.5 Limitations

We follow the categorization framework proposed by Runeson and Höst [3] to outline the limitations of the conducted evaluation, addressing potential threats to internal, external, and construct validity:

#### 5.6 Discussion

Overall, the direct automated feedback system was positively received by students, indicating its effectiveness in enhancing their programming skills, despite some areas for improvement.

### 6 CONCLUSION & FUTURE WORK

Future work includes enhancing the visualization of feedback, such as grouping and color coding the feedback items to make it easier to differentiate between critical feedback items, suggestions for improvement, and positive feedback. A high priority will be on further improving the overall quality of the feedback provided. We also aim to extend the implementation to support direct automated feedback for the remaining exercise types of Artemis. Another crucial step is to test direct automated feedback in a real-world setting by utilizing this feature in an actual course. This will allow us to collect comprehensive data to thoroughly evaluate the impact on student performance and motivation.

### REFERENCES

- [1] I Elaine Allen and Christopher A Seaman. 2007. Likert Scales and Data Analyses. *Quality progress* 40, 7 (2007), 64–65.
- [2] Marcelo Guerra Hahn, Silvia Margarita Baldiris Navarro, Luis De La Fuente Valentin, and Daniel Burgos. 2021. A Systematic Review of the Effects of Automatic Scoring and Automatic Feedback in Educational Settings. *IEEE Access* 9 (2021), 108190–108198. <https://doi.org/10.1109/ACCESS.2021.3100890>
- [3] Per Runeson and Martin Höst. 2009. Guidelines for Conducting and Reporting Case Study Research in Software Engineering. *Empirical Software Engineering* 14, 2 (April 2009), 131–164. <https://doi.org/10.1007/s10664-008-9102-8>
- [4] Roel J. Wieringa. 2014. *Design Science Methodology for Information Systems and Software Engineering*. Springer Berlin Heidelberg.