Research Article

# Development of a Basic Chemistry Conversational Corpus

**Maurice HT Ling[1]\*, Syameer Musttakim[1] and Poh Nguk Lau[1,2]**

[1]*School of Applied Sciences, Temasek Polytechnic, Singapore*
[2]*Learning Academy, Temasek Polytechnic, Singapore*

**\*Corresponding Author:** Maurice HT Ling, School of Applied Sciences, Temasek Polytechnic, Singapore.

## Abstract

Chatbot technology can be an important tool and supplement to education, leading to explorations in this area. Corpus-based chatbot building has a relatively low entry barrier as it only requires a relevant corpus to train a chatbot engine. The corpus is a set of human-readable questions and answers and may be an amalgamation of existing corpora. However, a suitable chemistry-based chatbot corpus catering for a freshman general chemistry course addressing inorganic and physical chemistry has not been developed. In this study, we present a basic chemistry conversational corpus consisting of 998 pairs of questions and answers, focused on a freshman general chemistry course addressing inorganic and physical chemistry. Ten human raters evaluated the responses of a chatbot trained on the corpus and suggests that the corpus resulted in better response than random (t = 17.4, p-value = 1.86E-53). However, only 20 of the 50 test questions show better responses compared to random (difference in mean score ≥ 1.9, paired t-test p-value ≤ 0.0324), suggesting that the corpus provides better responses to certain questions rather than overall better responses, with questions related to definitions and computational procedures answered more accurately. Hence, this provides a baseline for future corpora development.

**Keywords:** Chemistry; Conversational Corpus; ChatterBot

## Introduction

Chatbot can be defined as a computer program that mimics human conversation [1] and has its roots in Turing Test [2]. In recent years, chatbots have been used in a variety of fields [3]; such as commerce [4], healthcare [5], and education [6-9]; especially during COVID-19 pandemic as a means to deliver services [10-12]. Fonna and Widyantoro [13] found that chatbots incorporated into tutorials can supplement pharmacology education. Kovacek and Chow [14] presented a chatbot catering to radiation safety education. Several studies also pointed to favourable student learning outcomes when chatbot are deployed to teaching and learning. For example, Atmosukarto., *et al.* [15] showed that by incorporating chatbot to an online chemistry course, course completion rate can be improved as the chatbot can cater to student's on-demand query needs, as compared to the limited online presence of a human tutor. Chatbots could also be used as a preparatory resource to complement online courses to augment learners' readiness for high-stakes assessment [16]. The interactive affordance and immediate feedback features provided by chatbots also facilitate self-directed learning and self-evaluation [17].

Of the 6 paradigms of chatbot building classified by Luo., *et al.* [3]; namely, template-based, corpus-based, intent-based, recurrent neural network-based, reinforcement learning-based, and hybrid; corpus-based appears to be the easiest and most flexible. A chat corpus can be built from existing questions and answers pair and corpus-based paradigm will allow a new chatbot to be implemented by retraining it on a new conversational corpus [18,19], which had been successfully demonstrated [20-22]. Another important advantage is that a corpus can be incrementally built or amalgamation of existing corpora. Hence, a chatbot can be considered as a chatbot engine (which can be defined as a software component that accepts a natural human language input, processes, and respond with an output in a natural human language [23]) trained on a corpus. Therefore, the advancement of chatbot technology and use cases can be the parallel track of advancement of chatbot engine technology, and the advancement of corpus.

In this study, we present a basic chemistry conversational corpus of 998 pairs of questions and answers, focused on a tertiary-level freshman inorganic and physical chemistry course; as such a corpus has not been widely developed. Ten human raters were

used to evaluate the responses of a chatbot trained on the corpus and suggests that the corpus resulted in better response than random (t = 17.4, p-value = 1.86E-53). Hence, this provides a baseline for future corpora development.

## Materials and Methods
### Corpus development

A conversational chemistry corpus, with contents based on a freshman general chemistry course addressing inorganic and physical chemistry was built from existing frequently asked questions (FAQs). Senior year students and faculty members were asked to contribute questions and answers related to the course. These questions and answers were converted into YAML format (https://yaml.org/).

### Corpus testing

The corpus was used to train a chatbot, which was built on Chatterbot (https://chatterbot.readthedocs.io) [23] and known as ChemBot. An untrained chatbot, known as DumbBot, was used for comparison. A set of 50 questions were posed to both ChemBot and DumbBot and each of the replies were evaluated by 10 human raters for quality [24]. The 50 questions were sourced from the 10 human raters to represent questions that may be asked by freshmen. The overall mean scores, mean scores by raters, and mean scores by questions were evaluated as mean scores is the most common metric of academic performance [25,26].

### Statistical analysis

T-test assuming unequal variance was used to compare the mean scores between ChemBot and DumbBot. 1-way ANOVA was used to compare mean scores between the raters and between questions. Paired t-test was used to compare the mean scores, which were paired by question or by rater. A p-value of less than 0.05 was statistically significant.

## Results and Discussion
### Corpus improves chatbot responses

A total of 998 pairs of questions and answers were compiled as corpus, which was then evaluated by comparing a trained chatbot (ChemBot) against an untrained chatbot (DumbBot). Each of the 10 raters scored the replies from both ChemBot and DumbBot on 50 questions. Mean scores were obtained by aggregating scores of the 50 questions per rater across the 10 raters (giving a total of 500 scores). Table 1 presents the mean score and standard error rated on ChemBot and DumbBot.

The mean score of ChemBot is 3.75 with standard error of 0.158 while the mean score of DumbBot is 1.01 with standard error of 0.004, the difference of which is statistically significant using t-

|  | ChemBot | DumbBot |
|---|---|---|
| Mean score | 3.75 | 1.01 |
| Standard error | 0.158 | 0.004 |

**Table 1:** Descriptive statistics for answer accuracy for ChemBot and DumbBot.

test assuming unequal variances (t = 17.4, p-value = 1.86E-53). This suggests that the developed corpus improves the response of ChemBot when compared to DumbBot, which is supported by previous studies [19,21,27-29] and is the fundamental motivation for corpus development [20,30].

### Comparable scores by each ratter on chembot

Of the 10 raters, rater R8 is the only chemistry lecturer (last author) with the rest being students. The mean score of ChemBot (n = 50) by each rater ranged from 3.24 with standard error of 0.504 (by rater R8) to 4.20 with standard error of 0.481 (by rater R6). Although the scoring between rater R6 and R8 is significant on paired t-test as paired by questions (p-value = 9.73E-3); the correlation of scores given by rater R6 and R8 is 0.738, which is higher than that between rater R1 and R8 (see Figure 1). Interestingly, paired t-test suggests that the scores given by rater R1 and R8 is not significant (p-value = 0.144) despite having the lowest correlation. 2-samples t-test assuming unequal variances between rater R1 and R8 (p-value = 0.416) and between rater R6 and R8 (p-value 0.171) are not significant. This is supported by 1-way ANOVA suggesting no significance in the mean scores across all raters (Figure 2; F = 0.238, p-value = 0.989). For DumbBot, 1-way ANOVA also suggests no significance in the mean scores across all raters (Figure 2; F = 1.235, p-value = 0.271).
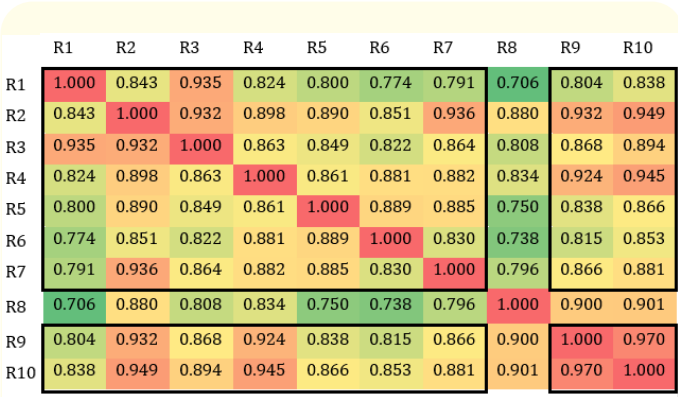


|  | R1 | R2 | R3 | R4 | R5 | R6 | R7 | R8 | R9 | R10 |
|---|---|---|---|---|---|---|---|---|---|---|
| R1 | 1.000 | 0.843 | 0.935 | 0.824 | 0.800 | 0.774 | 0.791 | 0.706 | 0.804 | 0.838 |
| R2 | 0.843 | 1.000 | 0.932 | 0.898 | 0.890 | 0.851 | 0.936 | 0.880 | 0.932 | 0.949 |
| R3 | 0.935 | 0.932 | 1.000 | 0.863 | 0.849 | 0.822 | 0.864 | 0.808 | 0.868 | 0.894 |
| R4 | 0.824 | 0.898 | 0.863 | 1.000 | 0.861 | 0.881 | 0.882 | 0.834 | 0.924 | 0.945 |
| R5 | 0.800 | 0.890 | 0.849 | 0.861 | 1.000 | 0.889 | 0.885 | 0.750 | 0.838 | 0.866 |
| R6 | 0.774 | 0.851 | 0.822 | 0.881 | 0.889 | 1.000 | 0.830 | 0.738 | 0.815 | 0.853 |
| R7 | 0.791 | 0.936 | 0.864 | 0.882 | 0.885 | 0.830 | 1.000 | 0.796 | 0.866 | 0.881 |
| R8 | 0.706 | 0.880 | 0.808 | 0.834 | 0.750 | 0.738 | 0.796 | 1.000 | 0.900 | 0.901 |
| R9 | 0.804 | 0.932 | 0.868 | 0.924 | 0.838 | 0.815 | 0.866 | 0.900 | 1.000 | 0.970 |
| R10 | 0.838 | 0.949 | 0.894 | 0.945 | 0.866 | 0.853 | 0.881 | 0.901 | 0.970 | 1.000 |

**Figure 1:** Correlation matrix between raters. The minimum Pearson's correlation coefficient is 0.706 (n = 50, t = 6.907, p-value = 1.021E-8) between R8 and R1. The minimum correlation among student raters (R8 is a lecturer) is 0.774 (n = 50, t = 8.469, p-value = 4.346E-11).
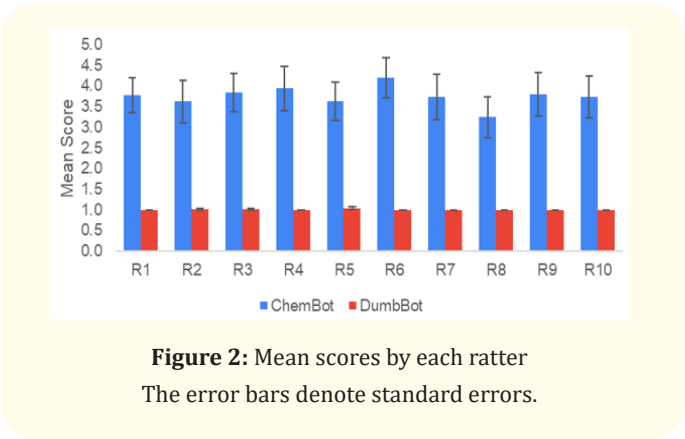
**Figure 2:** Mean scores by each ratter
The error bars denote standard errors.

After removing rater R8, the mean score of ChemBot (n = 50) by each student rater ranged from 3.62 with standard error of 0.523 (by rater R2) to 4.20 with standard error of 0.481 (by rater R6).

2-samples t-test assuming unequal variances between rater R2 and R6 is not significant (p-value = 0.416) with 1-way ANOVA suggesting no significance in the mean scores across all 9 student raters (Figure 2; F = 0.122, p-value = 0.998). For DumbBot, a 1-way ANOVA also suggests no significance in the mean scores across all raters (Figure 2; F = 1.200, p-value = 0.297). Taken together, the scores across all 10 raters are comparable.

**Mean scores differ by questions**

The mean score of ChemBot (n = 10) by question ranged from 1.00 with standard error of zero for 13 questions (Questions 7, 11, 18, 23, 24, 25, 27, 32, 33, 43, 44, 45, and 50) to 9.10 with standard error of 0.233 for question 26 (see Table 2 and Figures 3A to 3C). This is supported by 1-way ANOVA suggesting significant differences in the mean scores across questions (Figure 3; F = 59.939, p-value = 1.1E-166).

| Question | ChemBot Mean Score | ChemBot Median | DumbBot Mean Score |
|---|---|---|---|
| Q1. What is the electronic configuration of nitrogen? * | 5.50 (0.992) | 7.5 | 1.00 (0.000) |
| Q2. What are the forces between ammonia? | 1.70 (0.616) | 1.0 | 1.00 (0.000) |
| Q3. Why do Protons only exist in the nucleus? * | 3.60 (0.653) | 5.0 | 1.00 (0.000) |
| Q4. Why are ionic bonds the strongest bond? | 2.10 (0.737) | 1.0 | 1.00 (0.000) |
| Q5. What are orbitals? * | 3.90 (0.936) | 3.5 | 1.00 (0.000) |
| Q6. How to calculate pH of a solution? | 1.40 (0.267 | 1.0 | 1.00 (0.000) |
| Q7. What are the bonds present in a covalent bonding? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q8. What is hydrogen bonding? | 2.00 (0.516) | 1.0 | 1.00 (0.000) |
| Q9. Why do ionic compounds have a high boiling point? * | 8.50 (0.307) | 8.0 | 1.00 (0.000) |
| Q10. Why are some acids strong and some are weak acid? | 1.30 (0.213) | 1.0 | 1.00 (0.000) |
| Q11. Why water have hydrogen bonding? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q12. What is the classical method of naming compounds? | 1.60 (0.499) | 1.0 | 1.00 (0.000) |
| Q13. What is the equilibrium constant? * | 8.60 (0.600) | 9.0 | 1.00 (0.000) |
| Q14. What is ionic equilibrium? | 2.20 (0.727) | 1.0 | 1.00 (0.000) |
| Q15. How to calculate molarity? * | 8.70 (0.335) | 8.5 | 1.00 (0.000) |
| Q16. What are the different types of bonding forces? | 1.50 (0.342) | 1.0 | 1.00 (0.000) |
| Q17. What is a buffer solution? | 1.20 (0.133) | 1.0 | 1.00 (0.000) |
| Q18. How do I make a buffer solution? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q19. What is equilibrium? * | 8.40 (0.221) | 8.5 | 1.00 (0.000) |
| Q20. How to find dilution factor? * | 7.60 (0.542) | 8.0 | 1.00 (0.000) |
| Q21. How to calculate mol? | 1.20 (0.133) | 1.0 | 1.00 (0.000) |
| Q22. What is atomic mass? | 2.90 (0.752) | 2.0 | 1.00 (0.000) |
| Q23. What is atomic weight? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q24. What is molecular weight? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q25. What is formula weight? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q26. How to calculate percentage composition? * | 9.10 (0.233) | 9.0 | 1.00 (0.000) |
| Q27. What is the Avogadro's number? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q28. What is molar mass? * | 8.50 (0.543) | 9.0 | 1.00 (0.000) |
| Q29. How do I find the number of moles? * | 8.30 (0.633) | 9.0 | 1.00 (0.000) |

| | | | |
|---|---|---|---|
| Q30. What is stoichiometry? * | 8.00 (0.558) | 8.5 | 1.00 (0.000) |
| Q31. What is a limiting reactant? * | 8.70 (0.335) | 9.0 | 1.00 (0.000) |
| Q32. What is an excess reactant? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q33. How to find precent yield? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q34. What is a monoprotic acid? | 1.20 (0.133) | 1.0 | 1.00 (0.000) |
| Q35. What is amphoteric? | 1.40 (0.306) | 1.0 | 1.00 (0.000) |
| Q36. What is the unit of pH? | 1.50 (0.401) | 1.0 | 1.00 (0.000) |
| Q37. What do pH and pOH measure? | 1.10 (0.100) | 1.0 | 1.00 (0.000) |
| Q38. What happens when a strong acid and a base react? | 1.20 (0.200) | 1.0 | 1.00 (0.000) |
| Q39. What is a conjugate base and conjugate acid? | 1.20 (0.200) | 1.0 | 1.20 (0.133) |
| Q40. What is buffer capacity? * | 8.60 (0.306) | 8.5 | 1.00 (0.000) |
| Q41. What is equivalence point? * | 8.50 (0.269) | 8.5 | 1.00 (0.000) |
| Q42. What is chemical kinetics? * | 8.50 (0.500) | 9.0 | 1.00 (0.000) |
| Q43. What is activation energy? | 1.00 (0.000) | 1.0 | 1.20 (0.133) |
| Q44. What does a catalyst do? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q45. How does the temperature affect the reaction rate? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |
| Q46. What is a reversible reaction? * | 7.90 (0.458) | 8.0 | 1.00 (0.000) |
| Q47. What is a dynamic equilibrium? * | 7.30 (0.943) | 8.0 | 1.00 (0.000) |
| Q48. What is Le Chatelier's principle? * | 8.70 (0.300) | 9.0 | 1.00 (0.000) |
| Q49. Why does the equilibrium shifts? | 1.10 (0.100) | 1.0 | 1.00 (0.000) |
| Q50. What is a weak acid? | 1.00 (0.000) | 1.0 | 1.00 (0.000) |

**Table 2:** Mean score by question in ChemBot and DumbBot. Standard error in brackets.

*mean scores significantly higher in ChemBot than Dumbot (paired t-test difference in mean score ≥ 1.9, p-value ≤ 0.0324).

For DumbBot, only 2 questions (Questions 39, and 43) have a mean score of more than zero (mean score = 1.20 with standard error of 0.133, see Figure 3C). Interestingly, the mean score for question 43 is higher in DumbBot (mean score = 1.20) than Chem-Bot (mean score = 1.00); however, this difference is not significant (paired t-test p-value = 0.168).

Of the 50 questions, the mean scores of 20 questions (Questions 1, 3, 5, 9, 13, 15, 19, 20, 22, 26, 28, 29, 30, 31, 40, 41, 42, 46, 47, and 48) were significantly higher in ChemBot as compared to Dumb-Bot (difference in mean score ≥ 1.9, paired t-test p-value ≤ 0.0324). This suggests that ChemBot provides better responses to certain questions rather than overall better responses, which has been previously demonstrated [31]. This is plausible as Callejas-Rodríguez., *et al*. [32] had shown that chatbot personality may be generated from training corpus. In general, it is also noted that ChemBot handles closed-ended concepts such as direct definitions (Q40, 41 and 42) and those related to computational procedures (such as Q15, 26, 29) more accurately. Questions that are more diffused, demand more in-depth explanation (Q18, Q45), has "compare/contrast" requirements (Q10 and Q39) or more than one keyword (such as "electronic configuration" and "nitrogen" in Q1; "atomic" and

"mass" in Q22 and Q23) are typically poorly answered. Previous studies have demonstrated the impact of keyword recognition in the accuracy of question and answer generated by chatbot applications [17].

**Applications and future work**

This baseline corpus could be integrated into a chatbot application currently on trial within the institution. This application is a commercial solution purchased on a licence-basis and requires a direct feed of FAQs for training purposes. The chatbot is integrated with an institution-wide communications interface, the Microsoft Teams, on which learners could pose a question to the chatbot. Instead of directly providing an answer, the chatbot engine directs another learner in the class to provide the answer. The instructor could choose to rate the learner-provided answer, archive a good answer into the FAQ database for future use or provide a more accurate response from the original database. Two clear advantages are apparent. One, the learner actively contributes to the learning of their peers, and secondly, the corpus could also be enhanced. Such a chatbot could then be implemented and tested in a chemistry course.
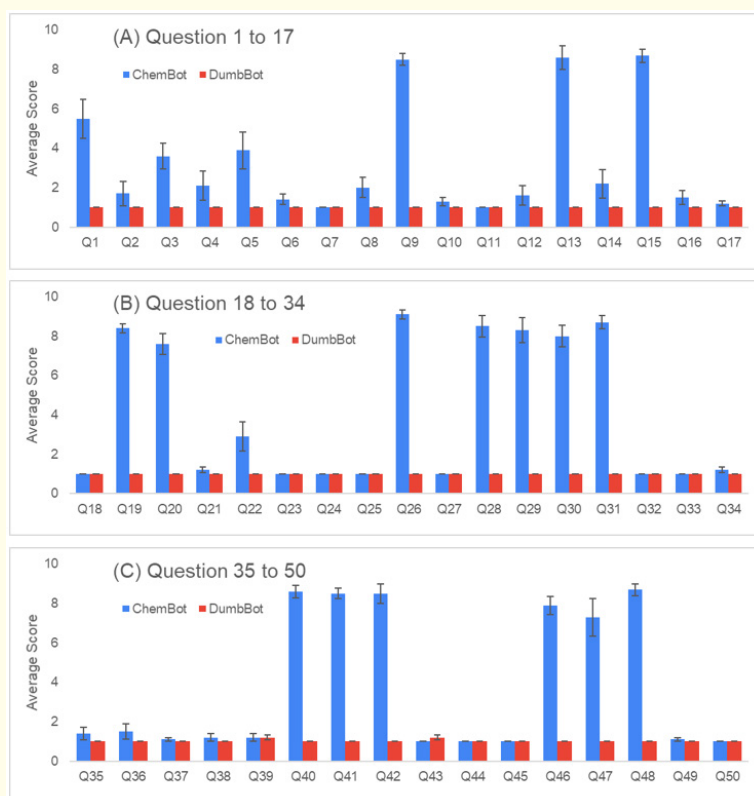
**Figure 3:** Mean scores across questions.

The error bars denote standard errors. Panel A shows mean scores from questions 1 to 17. Panel B shows mean scores from questions 18 to 34. Panel C shows mean scores from questions 35 to 50.

## Conclusion

Here, we present a basic chemistry conversational corpus consisting of 998 pairs of questions and answers, focused on tertiary year one inorganic and physical chemistry course; which was used to train chatbot based on Chatterbot engine and shown to generate better chatbot responses than untrained chatbot (t = 17.4, p-value = 1.86E-53). However, only 20 of the 50 test questions show better responses compared to random (paired t-test p-value ≤ 0.0324), suggesting that the corpus training results in better responses to certain questions rather than overall better responses. Hence, this study provides a baseline for future corpus development.

## Supplementary Materials

Materials from this study can be downloaded at https://bit.ly/ChemBot_1. Video showing comparative testing of ChemBot_1 and DumbBot can be found at https://youtu.be/tEJVRFphtLE.

## Conflict of Interest

The authors declare no conflict of interest.

## Bibliography

1. Adamopoulou E and Moussiades L. "An Overview of Chatbot Technology". AIAI 2020: Artificial Intelligence Applications and Innovations, eds Maglogiannis I, Iliadis L, Pimenidis E (Springer International Publishing, Cham) (2020): 373-383.

2. Turing AM. "Computing Machinery and Intelligence". *Mind* LIX 236 (1950): 433-460.

3. Luo B., *et al.* "A Critical Review of State-of-the-Art Chatbot Designs and Applications". *WIREs Data Mining and Knowledge Discovery* 12 (2022): e1434.

4. Landim ARDB. "Chatbot Design Approaches for Fashion E-commerce: An Interdisciplinary Review". *International Journal of Fashion Design, Technology and Education* 15.2 (2022): 200-210.

5. Tjiptomongsoguno ARW. "Medical Chatbot Techniques: A Review". Software Engineering Perspectives in Intelligent Systems, Advances in Intelligent Systems and Computing., eds Silhavy R, Silhavy P and Prokopova Z. "(Springer International Publishing, Cham) 1294 (2020): 346-356.

6. Okonkwo CW and Ade-Ibijola A. "Chatbots Applications in Education: A Systematic Review". *Computers and Education: Artificial Intelligence* 2 (2021): 100033.

7. Yang S and Evans C. "Opportunities and Challenges in Using AI Chatbots in Higher Education". Proceedings of the 2019 3rd International Conference on Education and E-Learning (ACM, Barcelona Spain) (2019): 79-83.

8. Hamam D. "The New Teacher Assistant: A Review of Chatbots' Use in Higher Education". HCI International 2021 - Posters, Communications in Computer and Information Science., eds Stephanidis C, Antona M, Ntoa S (Springer International Publishing, Cham) 1421 (2021): 59-63.

9. Akinwalere SN and Ivanov V. "Artificial Intelligence in Higher Education: Challenges and Opportunities". *Border Crossing* 12.1 (2022): 1-15.

10. Amiri P and Karahanna E. "Chatbot Use Cases in the Covid-19 Public Health Response. *Journal of the American Medical Informatics Association* 29.5 (2022): 1000-1010.

11. Zhu Y., *et al*. "It Is Me, Chatbot: Working to Address the CO-VID-19 Outbreak-Related Mental Health Issues in China. User Experience, Satisfaction, and Influencing Factors". *International Journal of Human-Computer Interaction* 38.12 (2022): 1182-1194.

12. Sweidan SZ., *et al*. "SIAAA-C: A Student Interactive Assistant Android Application with Chatbot During COVID-19 Pandemic". *Computer Applications in Engineering Education* 29.6 (2021): 1718–1742.

13. Fonna MR and Widyantoro DH. "Tutorial System in Learning Activities Through Machine Learning-Based Chatbot Applications in Pharmacology Education". 2021 8th International Conference on Advanced Informatics: Concepts, Theory and Applications (ICAICTA) (IEEE, Bandung, Indonesia) (2019): 1-6.

14. Kovacek D and Chow JCL. "An AI-Assisted Chatbot for Radiation Safety Education in Radiotherapy". *IOP SciNotes* 2.3 (2021): 034002.

15. Atmosukarto I., *et al*. "Enhancing Adaptive Online Chemistry Course with AI-Chatbot". 2021 IEEE International Conference on Engineering, Technology and Education (TALE) (2021): 838-843.

16. Korsakova E., *et al*. "Chemist Bot as a Helpful Personal Online Training Tool for the Final Chemistry Examination". *Journal of Chemical Education* 99.2 (2022): 1110-1117.

17. Mahroof A., *et al*. "An AI based Chatbot to Self-Learn and Self-Assess Performance in Ordinary Level Chemistry. 2020 2nd International Conference on Advancements in Computing (ICAC) (IEEE, Malabe, Sri Lanka) (2020): 216-221.

18. Shawar BAA. "A Corpus Based Approach to Generalise a Chatbot System". Doctor of Philosophy (University of Leeds, School of Computing) (2005).

19. Shawar BA and Atwell ES. "Using Corpora in Machine-Learning Chatbot Systems". *International Journal of Corpus Linguistics* 10.4 (2005): 489-516.

20. Rikters M., *et al*. "Designing the Business Conversation Corpus". Proceedings of the 6th Workshop on Asian Translation (Association for Computational Linguistics, Hong Kong, China), 54-61.

21. Shawar BA and Atwell E. "Using the Corpus of Spoken Afrikaans to Generate an Afrikaans Chatbot". *Southern African Linguistics and Applied Language Studies* 21.4 (2003): 283-294.

22. Shawar BA and Atwell E. "Arabic Question-Answering via Instance Based Learning from an FAQ Corpus". Proceedings of the CL 2009 International Conference on Corpus Linguistics". *UCREL* 386 (2016): 1-12.

23. Sim KS and Ling MH. "Installation and Documentation Evaluation of Recent (01 January 2020 to 15 February 2021) Chatbot Engines from Python Package Index (PyPI)". *Acta Scientific Computer Sciences* 3.8 (2011): 38-43.

24. Shawar BA and Atwell E. "Different measurements metrics to evaluate a chatbot system". (Association for Computational Linguistics) (2007): 89-96.

25. Kumar S., *et al*. "Defining and Measuring Academic Performance of Hei Students - A Critical Review". *Turkish Journal of Computer and Mathematics Education* 12.6 (2021): 3091-3105.

26. Alyahyan E and Düştegör D. "Predicting Academic Success in Higher Education: Literature Review and Best Practices". *International Journal of Educational Technology in Higher Education* 17.1 (2020): 3.

27. Kim J., *et al*. "Two-Step Training and Mixed Encoding-Decoding for Implementing a Generative Chatbot with a Small Dialogue Corpus. Proceedings of the Workshop on Intelligent Interactive Systems and Language Generation (2IS and NLG) (Association for Computational Linguistics, Tilburg, the Netherlands) (2018): 31-35.

28. Kowsher Md., *et al*. "Doly: Bengali Chatbot for Bengali Education". 2019 1st International Conference on Advances in Science, Engineering and Robotics Technology (ICASERT) (IEEE, Dhaka, Bangladesh) (2019): 19.

29. Blanc C., *et al*. "FlauBERT vs. CamemBERT: Understanding Patient's Answers by a French Medical Chatbot". *Artificial Intelligence in Medicine* 127 (2022): 102264.

30. Shawar BAA and Atwell E. "Automatic Extraction of Chatbot Training Data from Natural Dialogue Corpora" (2016): 29-38.

31. Kapočiūtė-Dzikienė J. "A Domain-Specific Generative Chatbot Trained from Little Data". *Applied Sciences* 10.7 (2020): 2221.

32. Callejas-Rodríguez Á., *et al*. "From Dialogue Corpora to Dialogue Systems: Generating a Chatbot with Teenager Personality for Preventing Cyber-Pedophilia. Text, Speech, and Dialogue, Lecture Notes in Computer Science., eds Sojka P, Horák A, Kopeček I, Pala K (Springer International Publishing, Cham) 9924 (2016): 531-539.