# USING LARGE LANGUAGE MODELS FOR ACADEMIC WRITING INSTRUCTION: CONCEPTUAL DESIGN AND EVALUATION OF THE SOCRAT PROJECT

Lukas Spirgi and Sabine Seufert
*University of St. Gallen*
*St.Jakob-Strasse 21, CH-9000 St.Gallen, Switzerland*

**ABSTRACT**

Academic writing has undergone significant evolution due to advancements in AI. Students are leveraging AI in diverse ways for their studies. This study introduces a course design (SOCRAT) to teach students genre-based academic writing through AI. Genre-based academic writing is an educational strategy instructing students in the writing techniques and norms pertinent to their specific academic disciplines. AI is utilised as a personal training system and research assistant in this proposed course design. Students require cognitive and metacognitive knowledge to effectively work with AI tools. The SOCRAT design is based on the concept of mastery learning to ensure that students build their competencies. An initial evaluation of the prompts developed for this design indicates that LLM can particularly assist students in analysing their written text and providing suggestions for enhancement. This can help students develop their genre-based writing skills. The models are not yet convincing for other types of tasks where the LLMs are required to give exact answers.

**KEYWORDS**

Academic Writing, Artificial Intelligence (AI), Text Genres, Large Language Model (LLM), Co-Creation

## 1. INTRODUCTION

Since November 2022, the swift proliferation of the robust chatbot ChatGPT has sparked widespread concerns, especially within academic institutions, regarding the potential misuse by students to produce texts without substantial effort. Potential drawbacks include the rise in plagiarism, surface-level learning, and reliance on such tools (Seufert et al., 2024). Several studies have examined how students benefit from ChatGPT and its impact on universities (Spirgi et al., 2024; von Garrel et al., 2023).

There are often calls for more oral examinations in the current discussion about using generative AI models such as ChatGPT in education. These are intended to ensure learners have analysed the topic in depth. Some universities even abolish qualification theses such as bachelor theses altogether (Zenthöfer, 2023).

Nevertheless, teaching academic writing remains essential in higher education, as it provides students with essential critical thinking and practical communication skills applicable across different fields. Incorporating large language models (LLMs) into academic writing pedagogy can enhance this pedagogical endeavour by providing tailored feedback and illustrative examples, thereby accelerating skill acquisition and proficiency (Seufert et al., 2024).

## 2. PROBLEM DEFINITION AND RESEARCH METHODOLOGY

As the introduction to the article shows, the advent of generative AI is massively transforming academic writing. Despite the support provided by AI tools, students must still be able to write academic texts. The article pursues the following goals based on the genre-based academic writing theory (Thoreau, 2006).

1. Developing a design for teaching students genre-based academic writing with the support of AI tools.

2. To assess how well current LLMs cope with the prompts developed for the design.

The research methodology utilised in this study is design-based research, known for its effectiveness in generating sustainable innovations in education and teaching (McKenney & Reeves, 2018). A design (SOCRAT) is being developed to show how academic writing skills can be taught in the first semester using AI in higher education. The design is based on the theory of genre-based writing and the skills students need to learn to work effectively with AI tools (Seufert et al., 2024; Thoreau, 2006). In order to identify the types of tasks that today's LLMs are good at, several prompts have been tested that can be used in the SOCRAT design. Criteria for evaluating the performance of each task were defined in advance. For the evaluation, 12 LLMs were selected. These models were selected based on performance, parameter size, and licence. The Hugging Face leaderboard was employed to evaluate the collective performance of the LLMs (Hugging Face, 2024). The interaction with the LLMs uses the zero-shot method, meaning the models must answer the question without context or examples. Table 1 lists the models tested. Each model is given a number. These numbers are used in the following.

Table 1. Selected LLMs

| Nr. | LLM | Nr | LLM | Nr. | LLM |
|---|---|---|---|---|---|
| 1 | Yi-34B-Chat | 5 | Claude 2.1 | 9 | OpenHermes-2.5-Mistral-7B |
| 2 | Mixtral-8x-7B-v0.1 | 6 | Openchat-3.5 | 10 | Nous-Hermes-13b |
| 3 | Mistral-Medium | 7 | Lama-2-70b-chat-hf | 11 | Gpt-4-1106-turbo-preview |
| 4 | Bard-jan-24-gemini-pro | 8 | Zephyr | 12 | Gpt-3.5-turbo-0125 |

Chapter 3 presents the theoretical foundations on which the SOCRAT design was developed. Chapter 4 presents the design developed and the results of the language models tested. A discussion of the implications follows this.

# 3. THEORETICAL BACKGROUND

## 3.1 Genre-Based Instruction to Develop Academic Writing Skills

Genre-based academic writing is a methodology in educational practice that focuses on teaching students the writing skills and conventions specific to their academic disciplines. This approach is precious in higher education, where developing research competence is crucial. The genre-based framework is grounded in the understanding that different academic fields have distinct forms of communication, and mastering these is key to academic and professional success (Seufert et al., 2024).

In higher education, students must learn various text genres to use correctly and complete their studies successfully. The term "genre" was coined by Thoreau (2006). He defines a genre as writing with a specific style, a specific target group of readers and a clear purpose. A genre, therefore, encompasses a functional perspective whereby specific social conventions, linguistic features and rhetorical structures of the text must be considered (Hyland, 2003). Each genre has its schematic structure. Genre approaches have been seen in recent decades as new ways of teaching and learning writing (Badger & White, 2000). They combine features and perspectives of more traditional approaches, such as the product-oriented (focus on linguistic knowledge, with particular attention to the appropriate use of vocabulary, syntax and cohesive devices) and the process-oriented approaches (focus on linguistic skills, e.g. drafting, planning, revising and editing text, rather than on linguistic knowledge). The genre approach has also become more prevalent in higher education in recent years to promote writing skills specifically (Kruse, 2016). Depending on the type of writing, students face varying expectations. However, students may struggle to fulfil these expectations due to the lack of standardised terminology for describing academic genres (Kruse et al., 2015).

A distinction can also be made between academic text genres with specific social functions, each of which entails different types of tasks in teaching (Seufert et al., 2024). Nesi and Gardner (2012) analyse various English universities' substantial assortment of student texts. They categorised these texts into 90 distinct genres, subsequently organised into 13 genre families. These genre families are based on five different social functions (Kruse, 2016):

- demonstration of knowledge and understanding (genres example: explanation)
- ability to make informed and independent arguments (genres example: essays)
- development of research skills (genres example: literature reviews)
- preparation for professional action (genres example: case reports)
- writing for oneself and others (genres example: narrative accounts)

An example of genre knowledge for writing an academic paper is the CARS model, which defines how a good introduction is written. John Swales formulated the CARS model after a comprehensive analysis of journal articles spanning diverse academic disciplines, aiming to enhance the quality of research introductions (J. M. Swales, 1990). The model outlines the typical structure used in crafting the introductory sections of scholarly research studies. According to the model, three rhetorical moves are employed to ascertain the research's background, motivation, and focal point: 1) Establishing a Territory, 2) Establishing a Niche, and 3) Occupying a Niche. For each of the three defined moves, the model shows how it can be implemented when writing an abstract (J. Swales, 2014)

The notion of text genres outlined here could serve as a foundational framework for elucidating and methodically fostering collaboration with AI-driven systems for research and writing within higher education (Seufert et al., 2024).

## 3.2 Knowledge Base for AI-based Academic Writing

To work effectively with AI writing tools, a variety of knowledge are required, as shown in Figure 1. (Seufert et al., 2024). The cognitive knowledge required consists of the following aspects (Seufert, 2024). Rhetorical knowledge pertains to the style's intent, the author's stance, and their understanding of the intended audience. Formal knowledge encompasses the culturally endorsed frameworks, etiquette, and norms regarding language selection within a given context. Subject-specific knowledge pertains to the information acquired within particular academic disciplines or areas of study. Procedural knowledge refers to how certain things are done, e.g. the steps involved in writing a text, such as researching and writing (Tardy et al., 2020). Interdisciplinary knowledge involves integrating concepts and content from various disciplines or subjects. Epistemic knowledge entails comprehending how seasoned professionals think and operate within their domains. This comprehension empowers learners to perceive the purpose and practicality of their acquired knowledge (OECD, 2020).

In addition to cognitive knowledge, metacognition is becoming increasingly crucial through collaboration with AI tools. Metacognition consists of two aspects: understanding one's thought processes and learning (metacognitive knowledge) and the skill to employ this understanding to manage and oversee one's learning journey (metacognitive regulation) (Seufert, 2024). Specifically, metacognitive knowledge is vital in assessing the outcomes produced by generative AI tools. Users must pose a series of inquiries to assess the output. The initial question concerns whether the output must adhere to truth (conditional knowledge). Subsequently, users must evaluate whether they possess the expertise required to determine the accuracy of the output (declarative knowledge) (Sabzalieva & Valentini, 2023).

Both students' cognitive knowledge and metacognition can be trained with AI tools. For example, an AI tutor can be used to train students' genre knowledge. At the same time, new technologies, incredibly generative AI with natural language dialogues, make it possible to promote better metacognitive knowledge, which was previously difficult to train. For example, learners can use generative AI as training systems, e.g. for a Socratic dialogue or as a dynamic evaluator to map thought processes and open up new ways of learning (Seufert, 2024).
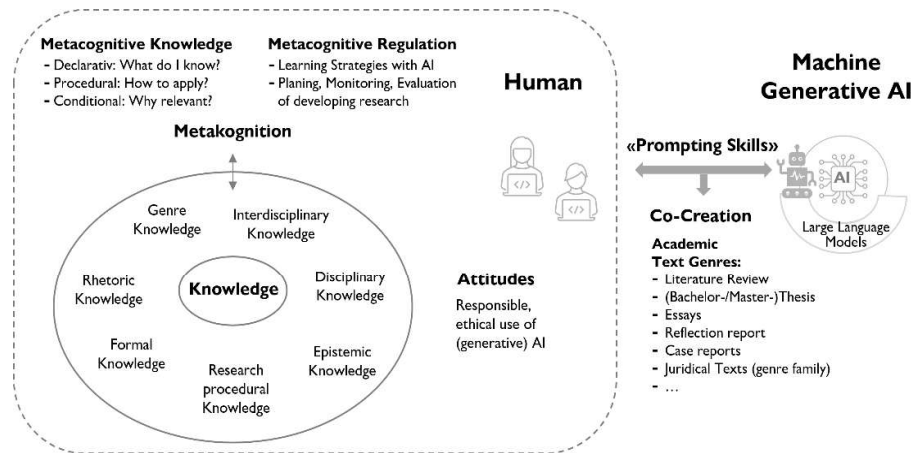
Figure 1. Type of knowledge for AI co-creation (Seufert, 2024)

# 4. RESULTS

## 4.1 Project SOCRAT (System of Critical Reasoning and Thesis)

### 4.1.1 Objectives of the SOCRAT Design

The SOCRAT (System of Critical Reasoning and Thesis) project aims to enhance students' research abilities during their initial semester at university. These essential skills are introduced through a compulsory introductory course attended by all students. In particular, the skills developed should make it easier for first-year students to start their studies. Additionally, the design empowers students to work with the AI tool.

### 4.1.2 SOCRAT as Training and Research Assistant

The course curriculum encompasses fundamental knowledge (as shown in Figure 2), such as various text genres and overarching research paradigms, as well as practical skills for designing, executing, and disseminating research projects. In the course's first part (Orientation Phase), students must build up the necessary specialist knowledge. This includes, above all, the cognitive knowledge presented in Chapter 3.1. In the second part of the course, students must apply what they have learnt by writing a systematic literature review. In this part of the course, students will go through the stages of planning their research, creating their insides and publishing their work.
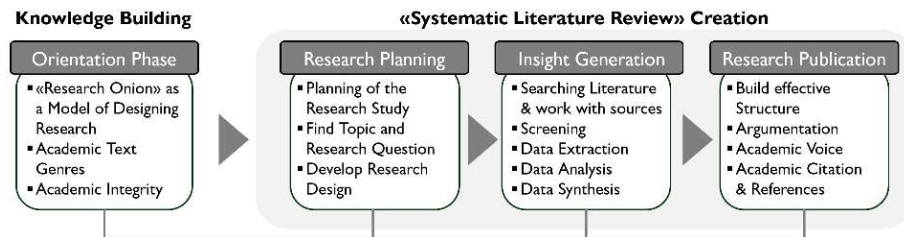


Figure 2. Stages and content of the SOCRAT design

The course consists of two parts and is structured according to Bloom's Mastery Learning (MT) concept. Mastery learning states that students should achieve a level of mastery in basic knowledge before moving on to new topics. If mastery is not achieved, students are given extra help to learn and revise, then retested (Bloom, 1968).

The use of AI in the two phases of the course is illustrated in Figure 3. First, the focus is on building research skills. In particular, the AI can provide personalised feedback, meaning each student is supported individually. The AI becomes a personalised training system. In addition to the necessary genre knowledge, students learn to use AI tools in academic writing. In particular, learning prompting skills is fundamental. At the end of the first phase, students sit a closed-book examination to test their acquired skills. Students who pass the exam enter the second phase of the course. The second part of the course requires students to apply and demonstrate these skills by working with AI to write their systematic literature review. The students write a scientific article in collaboration with the AI. This corresponds to the concept of "augmentation" (Jarrahi, 2018). This concept aims to create a working partnership between man and machine, where both can contribute their respective strengths. (Jarrahi, 2018; Seufert et al., 2019).
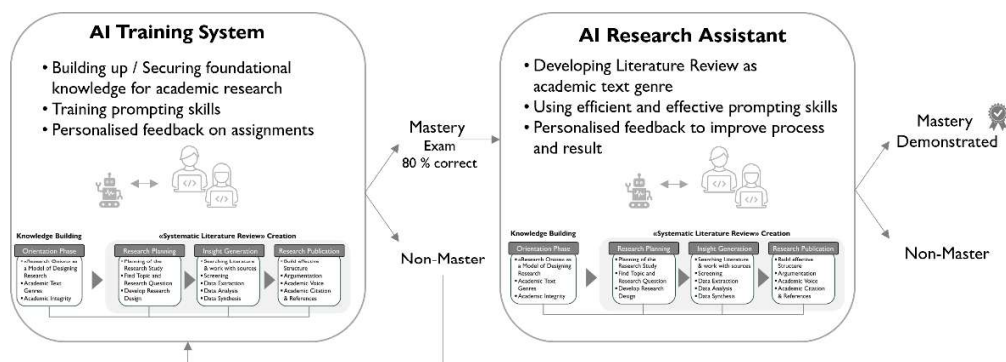
Figure 3. SOCRAT Design

The CARS model exemplifies the genre knowledge taught in SOCRAT Design. This model can illustrate an AI's role as a training system or a research assistant. Crafting an introduction for an academic article is a pivotal step in producing scholarly work. Structuring an introduction within a systematic literature review is an aspect of genre knowledge. When AI is utilised as a personal training system, it elucidates the principles of the CARS model and its potential applications within a systematic literature review framework. As a personal research assistant, AI answers students' queries about their texts and provides suggestions for more effective implementation of the CARS model in their systematic literature review introduction.

## 4.2 Evaluation of the Prompts Design Created for the SOCRAT Design

In the course design (SOCRAT) presented, AI plays an important role: in the first part as a personalised training system and in the second part as a personalised research assistant. The students use different prompts in the two roles of the AI. The next step is to test which types of prompts work well with the current state of the technology. The CARS model (see Chapter 3.1) is used as an example of genre knowledge to assess how well current LLMs cope with the prompts developed for the design. The performance of LLMs will be evaluated through two test series. The first trial contains tasks from lower taxonomy levels and is used to sort out unusable LLMs. The LLMs had to evaluate already written introductions in the second test series. The second test is exclusively assigned to language models who have completed the first task.

Table 2. Test procedure overview

| Test series | | Sample task for the LLMs |
|---|---|---|
| 1 | (9 tasks each in German and English) | • Name the 3 main moves in the CARS (Creating a Research Space) model by Swales. <br> • What does the acronym CARS stand for in Swales's CARS model? |
| 2 | (6 tasks each in German and English) | • The following text is an introduction to a research paper. Analyse this introduction according to the CARS model developed by Swales. For the three moves, "establish a territory", "establish a niche", and "occupying the niche", state exactly where they begin and where they end. For each move, name every single phrase that is typical for this move. |

### 4.2.1 Evaluation of LLMs in Test Series 1

As shown in previous chapters, this test requires LLMs to answer questions about the CARS model. The tasks with a low taxonomy level were sometimes solved very inaccurately. The LLMs are rated as adequate and inadequate based on their performance. A total of five models were rated unsatisfactory, and seven models were rated sufficient. Models 1, 6, 7, 8, 9, and 10 are unsuitable for the SOCRAT design due to poor performance in the first task. Specifically, Model 1 is imprecise, Model 6 has a high error rate, Model 7 struggles with German, and Models 8 and 9 answer too few questions correctly. Models 2, 3, 4, 5, 11, and 12 suit the SOCRAT design. Model 11 answered all questions correctly, while Model 12 also performed very well.

### 4.2.2 Evaluation of LLMs in Test Series 2

The LLMs were given a pre-written introduction, based on which the LLM had to solve the tasks described in Table 3. These are all tasks of higher taxonomy levels ('analysis' or 'evaluation'). In this section, the prompts used for the test are shown in italics.

Table 3. Tasks of the second test series

| Task set | |
|---|---|
| 1 | Identify and label the three CARS moves in a given instruction. |
| 2 | Identify which of the three moves is missing in the given instruction. |
| 3 | Complete the introduction based on input from the student. |
| 4 | Support in improving the introduction by giving Feedback |
| 5 | Compare introductions and improve one of them. |
| 6 | Find weak points in the introduction and improve them. |

The first task, in which the LLMs had to mark a CARS move's exact start and end points, was poorly solved. The LLMs could not achieve the accuracy promoted in the prompt: *"...For the 3 moves "establish a territory", "establish a niche", "occupying the niche", state exactly where they begin and where they end..."*. Task 2 was particularly badly solved when the prompts were written in English. All models have not understood the task. All LLMs solved task 3 very well. The prompt for task 3 described the task as follows: "*You are a tutor in a university course on academic writing. The students have learnt to write introductions according to the CARS model by Swales, but they need some help from time to time. The following is an incomplete introduction...*". The assistance provided by the LLMs was very valuable. The support of the LLMs was precious to the students. The LLMs wrote excellent supplements to the given introduction. The prompt for task 4 was built similarly to the prompts for task 3. Task 4 was also solved well by all LLMs. The prompt explicitly stated that no complete solution should be proposed. *"... Provide the student with suggestions about how she can edit the text. Only write a few sentences on your own. For the major part of the sentences, you only provide support for the editing, but no complete solution...".* The feedback provided is considered to be very helpful for the students. Tasks 5 and 6 were moderately well-solved. Many weaknesses were inherent in the instructions given to the LLM for completing tasks 5 and 6. All models offered valuable improvement suggestions. However, not all weaknesses were identified. It could be presumed that the LLMs would have discovered even more weaknesses if they had been queried again. Nevertheless, this was not pursued due to the zero-shot approach.

The following observations can be made on the individual models: Model 2 distinguished itself with coherent and concise answers in this evaluation, while Model 3 offered clear and actionable feedback. The lengthy responses of Model 4 were less disruptive, and Model 5 provided the most valuable feedback with its concise answers. Model 11 continued to be the top performer, whereas the feedback from Model 12 was perceived as less valuable.

## 5.  DISCUSSION AND OUTLOOK

The SOCRAT design was developed based on the theory of genre-based writing. Care has been taken to foster the necessary skills for students to work with AI tools. Design not only promotes cognitive knowledge but also metacognitive knowledge. Structuring the courses according to Bloom's mastery learning concept ensures all students achieve the required competencies.

The two test series demonstrate the technical capabilities of the LLMs in performing the tasks of a personal training system and a personal research assistant. The objective was to ascertain whether the SOCRAT design could be implemented. The findings indicate that, in principle, six different LLMs can execute the tasks. Interestingly, it appears that even free language models can accomplish the assigned tasks, which could be particularly beneficial for educational institutions given the often limited public budgets.

The language models had to solve tasks at different levels of Bloom's taxonomy (Bloom, 1976). LLMs better solve certain types of tasks than others. In particular, tasks 3 and 4 of the second test series were solved very well by the LLMs. The LLMs must provide feedback or write texts based on student input in these two tasks. In this type of task, the focus is less on the precision of the answers and more on whether the LLMs relate well to the given introduction. In these types of tasks, the output of the language models can vigorously promote the development of students' competencies by conveying cognitive knowledge and stimulating metacognition. The LLMs provide food for thought and encourage students to scrutinise their work.

The LLMs still have difficulties solving task types, such as task 1 of the current test series. The LLMs cannot mark the exact start and end of a move of the CARS model. Precision is difficult for the LLM to realise. This problem can also be seen in the first test series in which certain LLMs have already been sorted out. A high level of response accuracy is essential, especially for tasks with a low taxonomy level.

Our experiments with the LLMs also show that the design of the prompt is particularly critical to the quality of the output. The language in which the prompt is written also significantly influences the quality of the output. Most LLMs were trained with more English than German data. It is necessary to empower students to formulate precise and accurate prompts. These skills should be taught alongside the technical knowledge in the initial part of the course. For students to formulate precise prompts, they must correctly understand the basics from the first part of the course, such as genre knowledge.

The zero-shot approach was used to test the LLMs, meaning no follow-up questions were asked. However, the authors assume that follow-up questions could improve the performance of the LLMs in text analysis.

## 6.  CONCLUSION

Numerous studies suggest that the advent of AI fundamentally transforms academic writing (Boyd-Graber et al., 2023; Seufert et al., 2024; Spirgi et al., 2024). Students already use AI tools in their studies (von Garrel et al., 2023). Despite this, it remains crucial for students to master the art of scientific research. The rapid advancements in AI necessitate adaptations in university teaching programs. This article presents the theoretical foundations of genre-based writing, defines the competencies required to work with AI tools, and introduces the SOCRAT design. The two test series demonstrated that several LLMs can already perform the tasks of a personal training system or serve as personal research assistants. However, the performance of the LLMs varies for different types of tasks. The LLMs excel when tasked with making suggestions for improvement. It turns out that LLMs can be effectively used as tutors to promote the development of students' competencies. In other areas, the performance of LLMs is subpar.

Recognising that students can only effectively use AI tools if they possess the necessary knowledge to formulate precise prompts and evaluate the results is crucial. The design of the course and the tasks play a pivotal role in developing the student's skills. The next step involves testing the SOCRAT design with a select group of students.

# REFERENCES

Badger, R., & White, G. (2000). A process genre approach to teaching writing. *ELT Journal, 54*(2), 153–160. https://doi.org/10.1093/elt/54.2.153

Bloom, B. S. (1968). *Learning for Mastery: Instruction and Curriculum* (BR-6-2556). https://files.eric.ed.gov/fulltext/ED053419.pdf

Bloom, B. S. (Ed.). (1976). *Beltz-Studienbuch: Vol. 35. Taxonomie von Lernzielen im kognitiven Bereich* (5. Aufl., 17. - 21. Tsd). Beltz.

Boyd-Graber, J., Okazaki, N., & Rogers, A. (2023). *ACL 2023 policy on AI writing assistance*. https://2023.aclweb.org/blog/ACL-2023-policy/

Hugging Face. (2024). *Open LLM Leaderboard*. https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard

Hyland, K. (2003). *Second language writing*. *ACLS Humanities E-Book*. Cambridge University Press. https://doi.org/10.1017/CBO9780511667251

Jarrahi, M. H. (2018). Artificial intelligence and the future of work: Human-AI symbiosis in organisational decision making. *Business Horizons*, *61*(4), 577–586. https://doi.org/10.1016/j.bushor.2018.03.007

Kruse, O. (2016). Wissenschaftliches Schreiben forschungsorientiert unterrichten. In A. Hirsch-Weber & S. Scherer (Eds.), *Wissenschaftliches Schreiben in Natur- und Technikwissenschaften* (pp. 29–53). Springer Fachmedien Wiesbaden. https://doi.org/10.1007/978-3-658-12211-9_3

Kruse, O., Meyer, H., & Everke Buchanan, S. (2015). *Schreiben an der Universität Konstanz : Eine Befragung von Studierenden und Lehrenden.* https://doi.org/10.21256/zhaw-55

McKenney, S., & Reeves, T. C. (2018). *Conducting Educational Design Research*. Routledge. https://doi.org/10.4324/9781315105642

Nesi, H., & Gardner, S. (2012). *Genres across the disciplines: Student writing in higher education*. *The Cambridge applied linguistics series*. Cambridge University Press.

OECD. (2020). *Lernkompass 2030: OECD-Projekt Future of Education and Skills 2030 Rahmenkonzept des Lernens*. https://www.oecd.org/education/2030-project/contact/OECD_Lernkompass_2030.pdf

Sabzalieva, E., & Valentini, A. (2023). *Chatgpt and artificial intelligence in higher education: Quick start guide*. Organisation des Nations unies pour l'éducation, la science et la culture (UNESCO); Institut international de l'UNESCO pour l'enseignement supérieur en Amérique latine et dans les Caraïbes / International Institute for Higher Education in Latin America and the Caribbean (IESALC). https://eduq.info/xmlui/handle/11515/38828

Seufert, S. (2024). Zukunft Bildung: Auswirkungen generativer KI auf Bildungssysteme. In S. Seufert & S. Handschuh (Eds.), *Generative Künstliche Intelligenz: ChatGPT und Co für Bildung, Wirtschaft und Gesellschaft* (1. Auflage, pp. 139–164). Schäffer-Poeschel Verlag.

Seufert, S., Burkhard, M., Gubelmann, R., Niklaus, C., & Handschuh, S. (2024). Hochschulbildung: KI-basiertes Forschen und Schreiben. In S. Seufert & S. Handschuh (Eds.), *Generative Künstliche Intelligenz: ChatGPT und Co für Bildung, Wirtschaft und Gesellschaft* (1. Auflage, pp. 197–214). Schäffer-Poeschel Verlag.

Seufert, S., Guggemos, J., Meier, C., & Helfritz, K. (2019). *Augmentation. Personalentwicklung in der digitalen Transformation - Ergebnisse einer empirischen Studie*. Deutsche Gesellschaft für Personalführung. https://www.alexandria.unisg.ch/handle/20.500.14171/98945

Spirgi, L., Seufert, S., Delcker, J., & Heil, J. (2024). Student Perspectives on Ethical Academic Writing with ChatGPT: An Empirical Study in Higher Education. *Proceedings of the 16th International Conference on Computer Supported Education*(Volume 2), 179–186.

Swales, J. (2014). Create a research space (CARS) model of research introductions. *Writing About Writing: A College Reader*, 12–15.

Swales, J. M. (1990). *Genre analysis: English in academic and research settings*. Cambridge University Press.

Tardy, C. M., Sommer-Farias, B., & Gevers, J. (2020). Teaching and Researching Genre Knowledge: Toward an Enhanced Theoretical Framework. *Written Communication*, *37*(3), 287–321. https://doi.org/10.1177/0741088320916554

Thoreau, M. (2006). *Write on track: A guide to academic writing*. Pearson Education New Zealand.

von Garrel, J. von, Mayer, J., & Mühlfeld, M. (2023). *Künstliche Intelligenz im Studium  Eine quantitative Befragung von Studierenden zur Nutzung von ChatGPT & Co*. https://opus4.kobv.de/opus4-h-da/frontdoor/deliver/index/docId/395/file/befragung_ki-im-studium.pdf

Zenthöfer, J. (2023, December 1). Erste Uni schafft Bachelorarbeiten ab. *Frankfurter Allgemeine Zeitung,* 2023. https://www.faz.net/aktuell/karriere-hochschule/hoersaal/ki-und-plagiate-erste-uni-schafft-bachelorarbeiten-ab-19353621.html