

Large Language Models – Text Analysis + Text Generation in the Era of Generative AI

Seraina Fischer, Bigna Schmid and Annamària Sréter, *HSLU & FFHS*

Abstract— This paper explores large language models (LLMs) from three main perspectives. First, it delves into the technical foundations of LLMs, with a focus on natural language processing (NLP) and the implementation of low-rank adaptations for fine-tuning. These methods will be utilized in the second phase of our project to develop a specialized model for a specific use case. Second, the paper reviews prominent large language models, along with notable small language models, highlighting their differences and unique characteristics. Third, it profiles key contributors to the research and development of LLMs. Building on this research foundation, the paper then examines the application of LLMs across various industries, particularly within the DACH region, with an emphasis on Switzerland. We analyze their use in academia and industry, focusing on four key sectors: Law, Education, Customer Service, and Swiss German Context. We present state-of-the-art models such as Noxtua, the first European legal LLM; educa AI, an education-focused LLM from Germany; and Spitch, a Swiss conversational AI platform for customer service. This sets the groundwork for the second phase of our project, where we will fine-tune an LLM for a specific Swiss context.

Index Terms—Language model, natural language processing, education, language generation

1 INTRODUCTION

In recent years, Large Language Models (LLMs) have achieved significant technological advancements, making them a central topic in research and application across various industries. Despite their potential and wide applicability, there are challenges and developmental gaps. This paper aims to provide a comprehensive overview of the developments of LLMs by analyzing their technical background and foundational aspects. Particular attention is paid to the role of key figures in the development of these technologies, comparing their approaches and highlighting differences and similarities.

Furthermore, renowned language models will be identified and analyzed to better understand their application areas and capabilities. Four different areas of application within the DACH region will be investigated. A special focus is placed on the development and application of LLMs in the educational field, with a particular interest in adapting these technologies to the Swiss educational context. This includes investigating how LLMs can be specifically trained and developed as prototypes for creating educational content and enhancing learning experiences.

This paper includes the findings from an interview with an experienced individual active in the development of LLMs in the education sector. Namely the German LLM educa AI, whose approach will be notably discussed. The central research question is: How can Large Language Models be effectively adapted to the specific needs of the education sector in a German-speaking context and used to enhance learning processes?

LLMs excel in various linguistic and cognitive tasks (Chang et al., 2024). They generate fluent, precise text with high coherence and clarity, crucial for natural language generation. LLMs perform impressively in

language understanding tasks like sentiment analysis, text classification, and handling factual inputs, showcasing their depth of understanding. Their robust reasoning abilities encompass arithmetic, logical, temporal, and mathematical reasoning, allowing them to solve complex problems. LLMs' strong contextual comprehension ensures their responses are coherent and contextually appropriate, enhancing their use in dynamic conversations.

2 LITERATURE OVERVIEW

2.1 Architecture and Technical Foundation of LLMs

Natural Language Processing (NLP) tasks such as text summarization, generation, and prediction are driven by language models (Rizvi, 2019). Therefore, they play a crucial role in NLP. The field of NLP is divided into two primary disciplines: Natural Language Understanding (NLU) and Natural Language Generation (NLG) (Khurana et al., 2023). The primary aim of NLP is to understand, analyze, and manipulate language data using computational tools. It may be applied to various tasks and areas such as, for example, machine translation, text categorization, spam filtering, information extraction, summarization, dialogue systems, and medicine.

Khurana et al. (2023) highlight the interconnection between linguistics and computer science, demonstrating how NLP serves as a bridge between understanding human language and the mechanical process of generating it. NLU is concerned with enabling computers to interpret and analyze natural language, extracting meaning from text and speech. This involves a deep dive into the linguistic components that construct language, such as phonology, morphology, syntax, semantics, and pragmatics. NLU technologies identify intents (purposes)

and entities (elements) within user input, enabling differentiated responses to various expressions of the same intent. Tokenization and lemmatization break down texts into manageable units (tokens) and reduce them to their word stems, facilitating standardized processing. Subsequent sentence parsing assigns grammatical categories to words, and word vectors analyze their relationships to understand sentence structure.

NLG focuses on the synthesis of human-like text from structured data (Khurana et al., 2023). This process can be delineated into four sequential steps: identifying the text's objective, devising a strategic plan to achieve these objectives, evaluating the situational context, and finally, the execution or realization of these plans. To generate meaningful text, NLG systems employ various components including a speaker or a generator, and a mechanism for content selection, textual organization, and resource selection. NLG automates natural language creation, essential in applications like chatbots, weather reports, news generation et cetera (Kohne et al., 2020). NLG systems construct coherent sentences by integrating fixed text templates with dynamic information from external interfaces, like weather services, enabling chatbots to provide contextually relevant responses.

To ensure varied and natural interactions, NLG systems use multiple templates that can be randomized or combined (Kohne et al., 2020). Advanced NLG services produce highly realistic texts for specific fields such as finance, sports, and traffic, often indistinguishable from human-written texts. Providers like AX Semantics, textengine.io, and textOmatic offer scalable, cloud-based models for efficient text generation based on predefined scenarios, ensuring consistent and engaging user experiences.

An essential aspect of understanding and generating natural language lies in the creation and application of Language Models (LMs) (Rizvi, 2019). LMs predict the probability of a sequence of words to determine the most accurate representation. There are two main types of language models: statistical language models and neural language models.

Statistical language models rely on underlying statistical methods such as N-grams, Hidden Markov Models, and linguistic rules (Rizvi, 2019). However, these models face limitations due to the significant computational power required and the sparsity of N-grams, which can lead to zero probability outputs for unseen word combinations. Consequently, there has been a transition to neural language models that utilize deep learning techniques. A groundbreaking example of such an advanced neural technique is the transformer-based model GPT-2, which has demonstrated superior performance in various NLP tasks.

Recent advancements in abstractive text summarization (ATS) have focused on deep reinforcement learning (RL) and transfer learning (TL) methodologies (Alomari et al., 2022). ATS, crucial in NLP, aims to generate concise versions of texts through extractive and abstractive methods, with the latter offering superior readability and coherence but presenting more challenges. Deep neural

sequence-to-sequence models, especially the encoder-decoder architecture, have shown significant progress despite issues with long-term dependencies and generalization. RL approaches optimize beyond maximum likelihood by focusing on novel rewards and combining extractive and abstractive models, while pre-trained language models (PTLMs) like BERT and PEGASUS have revolutionized ATS with robust representations. Despite progress, challenges remain in improving factual accuracy, handling long documents, and enhancing novelty and diversity, with future research focusing on sophisticated models and hybrid methods for human-like summarization quality.

2.2 History of NLP

NLP research began in the late 1940s with machine translation (Khurana et al., 2023). Despite setbacks, it evolved with AI advancements like BASEBALL Q-A, LUNAR, and SHRDLU. The 1980s saw improvements in computational grammar and user intent processing, leading to systems like SRI's Core Language Engine. The 1990s focused on information extraction and automatic summarization, indicating practical applications.

In the early 2000s Bengio et al. introduced neural language modeling with innovations like feedforward neural networks and multitask learning (Bendig et al., 2001, as cited in Khurana et al., 2023). Advances by Mikolov et al. in word embedding and Sutskever et al. in sequence-to-sequence mapping further enhanced NLP capabilities (Mikolov et al., 2001; Sutskever et al., 2014, as cited in Khurana et al., 2023).

Transformers, like BERT, revolutionized NLP with attention mechanisms for context understanding, improving language modeling and machine translation (Zhao et al., 2023). Today, NLP tools span applications from sentiment analysis to semantic role labeling (Khurana et al., 2023). Recent advancements in NLP have significantly improved both the understanding and generation of natural language, incorporating diverse approaches and technologies. Noam Chomsky's rationalist and symbolic approaches, which suggest that certain knowledge is innate, have influenced the development of machines that mimic human cognitive processes by pre-installing fundamental linguistic rules. Statistical and machine learning models, including Naive Bayes classifiers and Hidden Markov Models (HMM), use algorithms to detect patterns and make predictions, facilitating tasks like document classification and speech recognition.

Since 2010, neural networks have revolutionized NLP. Word embeddings, Long-Short Term Memory (LSTM) networks, and Convolutional neural networks (CNNs) have enhanced language processing and generation, improving applications from text summarization to machine translation. The introduction of models like BERT, which analyzes text bidirectionally, has significantly improved context understanding (Khurana et al., 2023).

Transformer models, as discussed by Zhao et al. (2023), emphasize scalable architecture and extensive pre-training on vast text corpora, enhancing language nuances understanding and applicability across various

NLP tasks. LLMs demonstrate advanced capabilities in generating coherent, contextually relevant responses, advancing beyond traditional modeling techniques.

Ji et al. (2024) explored LLMs for generative recommendation (GenRec) systems, moving away from traditional discriminative methods that rely on designed user or item IDs. They propose using LLMs to generate recommendations directly from raw text data. Their GenRec LLM uses text-based item identifiers to generate personalized recommendations, leveraging semantic content for more accurate suggestions. Experimental results show GenRec outperforms traditional models like GRU4Rec and SASRec, especially with rich interaction data.

LoRA introduces low-rank matrix approximation into pre-trained models, significantly reducing computational load and storage requirements (Cabello, 2023). LoRA allows operating with lower resources while maintaining functionality. By freezing pretrained model weights and injecting trainable rank decomposition matrices, LoRA reduces trainable parameters for downstream tasks (Hu et al., 2021).

3 STATE-OF-THE-ART LARGE LANGUAGE MODELS

Recent advancements in text analysis and generation techniques using large language models (LLMs) have been remarkable (Khurana et al., 2023). Generative AI models like OpenAI's Chat GPT-4 and Meta's LLaMA have significantly enhanced NLP capabilities. This section provides an overview of state-of-the-art LLMs, including open-source models, smaller LMs, and LM developments specific to Switzerland.

3.1 Open-Source Models

Most state-of-the-art LLMs, like OpenAI's Chat GPT are massive, proprietary, closed-source models. This restricts cutting-edge LLM development to a few key players, prompting initiatives to make AI more accessible for public research and development to align advancements with human interests (IBM, n.d.).

Qwen 1.5 is collection of transformer-based models optimized for text generation and chat applications and was developed by Alibaba Cloud (Deci Research team, 2024). The model supports multiple language and is fine-tuned using Direct Preference Optimization and Proximal Policy Optimization and demonstrates robust performance across various benchmarks.

Abacus AI's Smaug series comprises models fine-tuned using DPO-Positive (DPOP), a variant of Direct Preference Optimization (Deci Research team, 2024). DPOP integrates a new term into the loss function to mitigate specific deficiencies in DPO. Smaug is notable for being the first open-source model to surpass an average score of 80% on the Open LLM Leaderboard, utilizing customized datasets during training to optimize performance on downstream tasks.

Introduced by the Hugging Face H4 Team, Zephyr 7B is a compact language model designed to align closely with user intent (Cabello, 2023). Despite its reduced size, it

performs superior to larger models, leveraging distilled Direct Preference Optimization (dDPO) and AI Feedback.

3.2 Small Language Models

The high costs and resource requirements for training large language models have led to the release of smaller, efficient models like Mistral and Orca (Huynh, 2024). These models achieve success through data quality and training efficiency.

Meta's LLaMA project, introduced in February 2023, aims to democratize access to powerful LLMs, making experimentation and innovation more accessible (Touvron, Lavril, et al., 2023; Touvron, Martin, et al., 2023). LLaMA models handle NLP tasks with less computational power, using public datasets to avoid proprietary data. LLaMA's models are trained on over 2 trillion tokens and are finetuned with reinforcement learning from human feedback.

Phi, developed by Microsoft Research, includes small transformer-based models trained on high-quality textbook data for tasks in NLP and coding (Microsoft, 2024). Phi demonstrates that emergent abilities in large models can also be achieved on a smaller scale through strategic training and data quality.

Orca, another Microsoft Research model, learns from GPT-4's rich signals and emulates the reasoning process, achieving parity with ChatGPT on various benchmarks (Mukherjee et al., 2023).

3.3 Swiss LLMs

In recent years, Switzerland has seen significant research and advancements in the development of LLMs (*SwissLLM*, n.d.). This progress reflects the country's commitment to advancing NLP technologies tailored to its unique linguistic landscape and industrial needs. These advancements reflect significant progress in LLM capabilities, making sophisticated NLP tools more accessible and efficient for diverse applications.

SwissLLM specializes in generative AI tailored to the needs of Swiss companies, ensuring privacy, security standards, and seamless data integration (*SwissLLM*, n.d.). They provide services like "creative agents" to support text and content creation for websites, blogs, and social media.

SwissBERT is a masked language model focused on Switzerland's four national languages (Vamvas et al., 2023). Trained on 21 million news articles, it excels in processing contemporary news articles and languages like Romansh Grischun, often underrepresented in other LLMs. It is based on the Cross-lingual Modular (X-MOD) transformer, enabling usage in 81 languages.

4 Key Individuals in the Evolution of LLMs

The development and advanced research in LLMs have redefined the boundaries of artificial intelligence, significantly transforming how machines process and generate language (Khurana et al., 2023). These remarkable technological advances stem not only from algorithms and datasets but predominantly from people. Researchers, developers, and visionaries have tirelessly worked and innovated to lay the foundations for these

revolutionary tools.

One of the most important figures in the development of large language models was Alan Mathison Turing, whose endeavors laid the foundation for today's LLM landscape (Gonçalves, 2023). In 1947, Turing declared in one of his lectures, "What we want is a machine that can learn from experience, and the possibility of letting the machine alter its own instructions provides the mechanism for this," underscoring the ability of LLMs to adapt, continuously enhance, and learn.

Machine Learning translation was one of the earliest uses of computers for language related tasks (Cabello, 2023). During World War II, Weaver and Booth began one of the first projects to use computers for language translation, providing a steppingstone for various research institutions.

The beginning of research in natural language processing dates back to the 1950s when researchers at IBM and Georgetown University worked on developing a system to automatically translate phrases from Russian to English (Cabello, 2023). This project marked one of the earliest ventures into machine learning language translation.

Until the 1960s, many researchers struggled to teach machines the intricacies of human language. However, with the introduction of ELIZA by MIT researcher Joseph Weizenbaum, the world saw the first chatbot (Rusch, 2023). This model used pattern recognition to simulate conversations, transforming user input into questions and generating responses based on predefined rules.

Weizenbaum introduced ELIZA as the world's first artificial psychiatrist (Tarnoff, 2023). With ELIZA Joseph Weizenbaum illuminated a mechanism of the human mind that strongly affects how we relate to computers. He used the concept of transference which describes human's tendencies to project feelings about someone from their past on to someone in their present to explain the success of ELIZA.

In 1970, Terry Winograd at MIT created SHRDLU, software capable of engaging in conversations about a confined virtual environment called the "blocks world," where users could interact by moving objects and posing questions (Cabello, 2023). The introduction of the Hidden Markov Model by Leonard E. Baum (1971) and Conditional Random Fields (CRFs) marked a shift from using strict rules to statistical models that analyze real-world text examples.

In 1996, Google made their PageRank algorithm public, which used N-gram models to estimate the likelihood of a word appearing by analyzing preceding words in a sentence (Cabello, 2023). Another significant milestone was the development of Recurrent Neural Network Language Models (RNNLM) by Elman in 1991, which could capture sequential dependencies in language. A few years later Landauer and Dumais leveraged high-dimensional semantic space to uncover concealed relations and meanings within textual data and thereby strongly promoted Latent Semantic Analysis.

In 1997, the German researchers Hochreiter and Schmidhuber introduced LSTM, which could solve

complex, artificial long time lag tasks that had never been solved by prior models (Hochreiter & Schmidhuber, 1997).

Bengio et al. introduced the first neural language model using a one-hidden-layer feed-forward neural network in 2003, pioneering word embedding (Cabello, 2023). Shortly after, Tomas Mikolov and his team at Google introduced Word2Vec.

Bahdanau and his colleagues from the institute Mila introduced sequence-to-sequence models in 2015, efficiently mapping variable-length input sequences to variable-length output sequences (Bahdanau et al., 2016).

One of the most critical steps in the development of LLMs was the publication "Attention Is All You Need," by Vaswani initiating the transformer era (Vaswani et al., 2023).

With the introduction of GPT-1 in 2018, OpenAI set a significant milestone for LLM development (Slats, n.d.). Shortly after, Google developed BERT (Bidirectional Encoder Representations from Transformers), highlighting the potential of pre-trained models. Nicolas Sornin, Olivier Sellier, and François Sforza developed LoRA in 2019, a groundbreaking mechanism for further LLM development.

Among the current luminaries in LLM development, the scientific community owes much to figures like Yoshua Bengio, Geoffrey Hinton, and Yann LeCun (LeCun et al., 2015). Their pioneering work in deep learning has laid the groundwork for modern natural language processing. Their research has not only shaped the architectures and training methods of today's LLMs but has also set the stage for future innovations.

Jeffrey Dean and Sanjay Ghemawat at Google have also made substantial contributions through their development of TensorFlow, a platform that has revolutionized how researchers and developers train and implement complex models, greatly enhancing the accessibility and scalability of LLMs (Abadi et al., 2018).

In the industrial sector, executives and technologists from companies such as OpenAI, Google, and Meta have significantly contributed to the development and popularization of LLMs (Brown et al., 2020). OpenAI, under the leadership of Sam Altman and Ilya Sutskever, has redefined text generation standards with its GPT models, especially GPT-3 and GPT-4. Their work demonstrates the potential of these technologies and pushes the boundaries of the conceivable.

At Google, the team led by Jeff Dean and Rajat Monga developed the BERT model, marking a substantial advancement in text data context understanding (Devlin et al., 2019). BERT has inspired numerous subsequent developments in LLMs and is extensively applied across various sectors.

Meta (formerly Facebook) has also contributed innovative research, particularly through the efforts of Joelle Pineau and Yann LeCun (Kraus et al., 2022). Their focus on self-supervised learning has facilitated significant progress in creating efficient and powerful LLMs.

5 Challenges and Limitations

Despite their tremendous potential, LLMs face several obstacles in everyday usage and adoption. This section

outlines the current challenges and limitations of LLMs.

Hallucinated Output and Spread of False Information

LLMs are known to generate "hallucinated" responses, which do not adhere to factual accuracy (Mündler et al., 2024). This issue becomes critically problematic when LLMs are deployed in sensitive areas where precision is essential. Mündler et al. (2024) highlight the risks associated with misinformation propagated by LLMs, underscoring the necessity for robust validation processes for the content they generate. Their study found that approximately 17.7% of sentences produced by LLMs, such as ChatGPT, contain contradictions, illustrating significant challenges in ensuring factual accuracy. They propose a prompting-based framework aimed at effectively detecting and mitigating these contradictions, significantly enhancing the overall reliability of LLM-generated text.

Kreps et al. (2020) explored the credibility and impact of AI-generated news in misinformation campaigns. Through three experiments, they assessed whether AI-generated texts are perceived as credible and their influence on public opinion regarding foreign policy. Their findings indicate that AI-generated texts curated by human editors were perceived as credible as real news articles. They found that while larger AI models produced more credible texts, improvements diminished with increasing model size.

Partisanship also significantly influenced perceived credibility, with individuals finding politically congenial stories more credible (Kreps, 2020). Disclaimers had mixed effects on perceived credibility. The implications of these findings are profound, as AI tools can generate credible-sounding misinformation at scale, posing a threat to democratic institutions by undermining trust in the media. Digital media literacy interventions and technological solutions, such as AI models detecting other AI-generated content, are recommended by the authors to combat this issue.

Ethical Considerations and Academic Integrity

LLMs in education pose ethical dilemmas, particularly around academic integrity. LLMs may encourage academic dishonesty if students use AI-generated content without proper citation (Agarwal et al., 2023). Institutions are advised to create clear policies for AI use in academic work, requiring explicit acknowledgment of AI assistance and defining strict consequences for violations (Yan et al., 2024).

Mollaki (2024) emphasizes integrating ethical considerations into AI deployment in education to maintain academic integrity. Accurate detection of AI-generated content is crucial to ensure submissions reflect students' knowledge. The European Commission (2022) outlines ethical guidelines for AI use in education, warning of severe penalties for non-compliance.

Benjamin Ledel, founder of Digital Learning GmbH and LLM educa AI, noted ethical concerns with their text-to-image converter, Image AI (B. Ledel, personal communication, 23 May 2024). Despite efforts to filter

content, misuse remains a concern due to diverse training data. educa AI restricts image generator use to teachers and has developed custom models for students with cognitive impairments to ensure a safer learning environment.

Bias and Fairness

LLMs can replicate and amplify biases present in their training data, leading to discriminatory outcomes in critical decision-making areas like recruitment and law enforcement. Bender et al. (2021) call for diverse data sets and bias-mitigation algorithms to promote fairness and equity in model outcomes.

Transparency and Explicability

The complexity of LLM algorithms leads to opaque decision-making processes. Burrell (2016) identifies three types of opacity: intentional secrecy, technical illiteracy, and intrinsic algorithmic opacity. They advocate for "explainable AI" to ensure stakeholders can understand, trust, and govern these technologies.

Privacy Concerns

LLMs' ability to process and store vast amounts of personal data raises privacy concerns. Shokri et al. (2019) highlight the risks of adversarial attacks that exploit model vulnerabilities to access private data. Robust security measures and privacy-preserving algorithms are essential for responsible data management.

Societal Impact of LLMs

Automation and the Job Market

LLMs' automation capabilities threaten jobs, especially in sectors relying on cognitive tasks. Arntz et al. (2016) discuss the dual impact of automation: enhancing economic efficiency while risking job displacement and income inequality. They recommend policies supporting education and skills training to help workers adapt to new roles.

Spread of Disinformation

LLMs are used to create and disseminate disinformation, influencing public opinion and electoral outcomes. Bradshaw & Howard (2019) document the global scale of social media manipulation and call for international cooperation to regulate AI technologies and develop countermeasures against their misuse.

6 LLMs in Education

Recent advancements in LLMs like GPT-3.5 and GPT-4 have revealed substantial untapped potential in various sectors, including education (Kasneci et al., 2023). These models exhibit versatility and adaptability that extend beyond corporate use, offering significant benefits across educational domains. While tools like Google Translate, Google Scholar, and Grammarly have long been integral to professional and academic processes, the emergence of ChatGPT has marked a significant shift in digital engagement due to its user-friendly interface and high-quality results (Gimpel et al., 2023). Despite its advantages, the use of ChatGPT presents challenges such as data transparency issues and the potential for

generating nonsensical outputs. These concerns highlight the need for careful integration of LLMs in educational settings.

BaiDoo-Anu & Owusu Ansah (2023) discuss the transformative potential of ChatGPT in education, noting its benefits in personalized tutoring, automated essay grading, and accessibility improvements. However, they also emphasize limitations like the lack of human interaction and potential biases in AI outputs. The effective integration of ChatGPT requires educators to adapt and employ these tools ethically and efficiently.

Kasneci et al. (2023) identify six main target groups for the usage of LLMs in education, focusing mostly on the learner's perspective. The target groups include primary and secondary students, university students, remote learners, learners with disabilities, and professional learners. Each group benefits from tailored applications, enhancing the efficiency and effectiveness of educational processes.

Javid et al. (2023) address the teacher perspective. They particularly explore the role of ChatGPT in creating personalized lesson plans and supporting academic tasks. They highlight concerns about plagiarism and the potential inaccuracies in AI-generated content, stressing the need for critical assessment of AI outputs by students.

Gan et al. (2023) consolidate research on LLMs in education, highlighting their potential in personalized learning, instructional support, and educational content creation. They also discuss challenges like privacy protection, data bias, and algorithm transparency, advocating for evolving teacher roles to adapt to these technologies.

Guidelines for the effective use of LLMs in education emphasize the importance of critical reflection and ethical use, recommending interactive dialogues with AI for customized educational support and stressing the necessity for rigorous verification of AI-generated information (Gimpel et al., 2023). While LLMs like ChatGPT offer substantial benefits for education, their integration requires careful consideration of their limitations and ethical implications. Collaboration among educators, policymakers, and technology experts is essential to develop strategies that maximize the advantages of LLMs while addressing their challenges.

LLM in Education in DACH Region

educa AI

In recent years, AI integration into educational settings has addressed various challenges. One notable innovation is educa AI, developed by Digital Learning GmbH in Germany (educa AI, 2021). Benjamin Ledel, the company's owner, researched and worked on leveraging AI long before the rise of ChatGPT and other LLMs. educa AI, designed primarily for students, educators, and professionals, utilizes authentic materials from private schools to ensure content integrity and relevance (Interviewee B. Ledel, personal communication, 23 May 2024).

educa AI adheres to strict data protection regulations, operating mainly in German. It serves as a support

database, a medical database, and a resource for curating and disseminating learning materials. Its functionalities include generating letters, reports, knowledge queries, and offering interactive features like memory exercises, multiple-choice questions, and transcription services in over 56 languages (educa AI, 2021). The platform includes AI-powered tools such as chat, image generation, audio generation, document generation, and a comprehensive knowledge database. However, the image generator, Image AI, only operates in English, posing a potential language barrier.

educa AI aims to address staff shortages and digital resource deficiencies while fostering a conducive learning environment (educa AI, 2021). It offers 24/7 virtual teacher availability, enhancing learning outside the classroom. The platform helps teachers improve their digital skills, enabling easy creation of quizzes and educational materials (Interviewee B. Ledel, personal communication, 23 May 2024).

educa AI differentiates itself by ensuring transparent tracking of data sources and prioritizing user privacy, through storing requests only with explicit user consent (educa AI, 2021). The model was initially trained on general data sources like Wikipedia and is finetuned with the client's data to ensure context-specific, trackable, and relevant output.

In its early stages, educa AI encountered challenges, such as addressing the formal and informal forms of "you" in German (Interviewee B. Ledel, personal communication, 23 May 2024). Through extensive manual labor and significant time commitment, these challenges were successfully addressed. Currently, educa AI measures its success through qualitative feedback from teachers and students, which has been positive so far. A planned study aims to compare educa AI with the open-source model ChatGPT, hoping to provide quantitative evidence of its performance.

Synte

The IU International University in Germany developed Synte, an AI-driven teaching and learning assistant aimed at personalizing education and improving learning outcomes (IU Internationale Hochschule, n.d.-a). Synte features a question-answer tool, exam preparation support, deep dialogue learning, and pre-assessment capabilities. It provides 24/7 access to a vast repository of course materials and promotes critical thinking and self-assessment (IU Internationale Hochschule, n.d.-b).

Accessible via web application, chat clients, and mobile devices to all IU students, Synte operates on state-of-the-art NLP models, utilizing GPT Version 3.5 for its technical infrastructure (IU Internationale Hochschule, n.d.-a). User feedback has been overwhelmingly positive, with 85% of users finding the platform's answers helpful, indicative of its growing acceptance and utilization relative to other educational channels. Unlike educa AI, Synte is not available for public use and can hence not be tested by non-students of IU.

Swiss Universities

The University of Zurich (UZH) offers resources on

using LLMs in education but has no proprietary LLM initiative (Universität Zürich, n.d.). Instead, educators and students use publicly available generative AI tools. UZH formed the learning laboratory "LeLa" with other Zurich-based universities to enhance digital literacy, conducting webinars on digital and generative AI tools in academics. UZH also provides links to the AI Campus, a platform dedicated to AI and data skills training.

ETH Zurich is exploring the development of a bespoke LLM tailored with unique educational materials in response to growing student demands (Walther, 2024). Students use public LLMs for tasks ranging from programming to creating study flashcards. Recognizing limitations in handling complex mathematical concepts and specialized subjects, ETH aims to address these gaps.

The SWISS AI Initiative, led by ETH Zurich and EPFL, involves multiple Swiss academic and research institutions (ETH AI Center & EPFL AI Center, n.d.). Supported by the supercomputer ALPS, this initiative focuses on AI research and developing large-scale AI systems across domains like science, education, health, and environmental sustainability. Research areas include model development, LLM security, privacy, scaling infrastructure, and advanced LLM strategies.

7 LLMs in the Legal Context

The widespread availability of legal documents in Switzerland, such as court decisions, laws, articles, commentaries, and contracts, provides a robust foundation for applying NLP to legal tasks (Niklaus et al., 2021). Legal NLP involves using NLP techniques on legal texts to streamline and enhance legal processes (Berner Fachhochschule, n.d.-a). The precise nature of legal language and the sensitivity of legal matters make it ideal for implementing language models to improve legal efficiency. Law firms, police, public prosecution offices, and courts can use LLMs to analyze contracts or legal jurisdictions, offering context-based analysis beyond simple keyword searches. BFH's research with the federal prosecution office shows that smaller, bespoke LLMs often outperform larger public models for legal queries due to their ability to be trained on specific internal data without compromising confidentiality (Berner Fachhochschule, n.d.-a).

Niklaus et al. (2021) explored Legal Judgment Prediction (LJP) to forecast court decisions based on facts and historical data, marking the first attempt to develop LJP models in German, French, and Italian. This initiative aimed to assist legal professionals in preparing arguments and prioritizing cases, thus enhancing efficiency and reducing legal case backlogs. The model, trained on data from 85,000 cases from the Federal Supreme Court of Switzerland (2000-2020), revealed performance issues with longer texts and across different law types, but no significant difference related to case chronology. Initially, a BERT model limited to 512 tokens was used, which was inadequate for the average case length. This was addressed by adapting to Long BERT, supporting up to 2048 tokens, albeit with increased processing demands. Poor performance in Italian, due to a limited dataset,

suggested benefits from using Adapters or cross-lingual transfer approaches to enhance data availability.

BFH (n.d.-b) outlines a three-sided approach for effectively using LLMs in legal settings: gathering high-quality data, ensuring legal compliance with data protection laws, and anonymizing data to mitigate the risk of replication by LLMs. The university also provides access to numerous datasets and pretrained models based on BERT architecture, which offer functionalities similar to AI models like ChatGPT (Berner Fachhochschule, n.d.-b).

Noxtua

Germany's largest commercial law firm, CMS, in collaboration with AI company Xayn, developed Noxtua, Europe's first sovereign legal AI model (Xayn, n.d.). Noxtua is tailored to European legal queries, assisting lawyers and legal professionals by analyzing, reviewing, summarizing, and generating text within legal documents, with a focus on European legislation such as the GDPR. Although primarily focusing on German and English, Noxtua supports additional languages and is trained on high-quality legal datasets labeled by law professionals, adhering to stringent confidentiality and security standards suitable for commercial use.

8 LLMs in a Swiss German Context

Swiss German, spoken by about five million people in Switzerland primarily in informal settings, comprises numerous dialects that lack standardization and differ significantly from Standard German (Plüss et al., 2020). These variations present challenges and opportunities for AI models, necessitating tailored approaches in language processing technologies to effectively handle these regional and diverse linguistic forms.

Few researchers have explored the potential of utilizing LLMs for Swiss German content due to the absence of a standardized written form (Gerlach et al., 2023). Gerlach et al. (2023) used a multilingual pre-trained model to translate automated speech recognition of Swiss German into Standard German, employing various segmentation methods from Swiss German TV shows. Despite progress, the model still fails to produce entirely accurate Standard German.

Honnet et al. (2017) developed a machine translation system for Swiss German, gathering over 60,000 written words and using a combination of neural and statistical machine translation methods. Performance decreased as test data diverged from training data due to the high variability of Swiss German influenced by geography and topics.

Köchli et al. (2023) attempted to create an LLM for Swiss German using a BERT model but faced challenges due to processing time, computational power, and the lack of a pre-trained model for Swiss German. They suggested improving accuracy by categorizing the language database by dialect and pre-training the model separately for each.

Thommen (2023) examined voice recognition models in Swiss German dialects, identifying word boundary identification and limited speech data availability as

primary challenges.

Industry & Public Sector

The Institute of Data Science at the University of Applied Sciences and Arts of Northwestern Switzerland (FHNW) is developing an Automatic Speech Recognition (ASR) model for Swiss German dialects, converting them into Standard German text (University of Applied Sciences and Arts of Northwestern Switzerland, n.d.). They compiled 293 hours of speech recordings but need thousands more hours to match established models. They introduced a multi-stage sentence reordering approach and a novel Intersection over Union (IoU) estimator for model refinement, making the dataset and model publicly available.

Recent high-quality datasets like SwissDial, SPC, or SDS-200 have the potential to advance Swiss German speech translation (Paonessa et al., 2023).

Whisper

Whisper, an open-source neural network for English speech recognition, has been adapted for Swiss German dialects (Radford et al., 2023). Swisscom, Switzerland's largest telecommunication provider, is developing a Swiss German speech system, focusing on Interactive Voice Response (IVR) and collaborating with IDIAP researchers for technical implementation (Widmer, 2018). They aim to handle diverse dialects by training with 3,000 hours of transcribed speech.

The Lucerne University of Applied Sciences used Whisper to extract spoken language from audio and video files in Swiss dialects, noting significant text understanding despite room for improvement (Hochschule Luzern – Design Film Kunst., n.d.).

Paonessa et al. (2023) utilized models XLS-R, Trafo, and Whisper to construct a Swiss German translation system. They found that dialects similar to others improved performance, while distinct dialects required in-dialect data. They highlighted the challenges of dialect transfer and managing differences between Swiss German and Standard German.

Spitch

Spitch, a global provider of Conversational AI platforms, offers a range of voice and text solutions, including virtual assistants, speech analytics, chat platforms, voice biometrics, and knowledge bases (Spitch, n.d.). They specialize in intent recognition using advanced NLP/NLU and sentiment analysis techniques, serving clients like Migros Bank to streamline customer verification processes (Spitch, 2020).

9 LLMs in the Customer Service Industry

The primary role of a service desk is to handle service requests for incident resolution, acting as a central contact point for organizational inquiries or as an interface between customers and service providers. Responsibilities include managing incidents, problems, configurations, transformations, and releases, directly communicating with users, employees, or customers (Steinig, 2023). Key competencies for service desk employees involve

categorizing service requests based on context and priority and understanding and resolving them. Customers expect quick, thorough answers and 24/7 service across multiple channels (Giovis & Rozsa, 2023).

Currently, AI is mainly used for categorization, relying on rule-based systems or traditional machine learning algorithms to automate tasks and provide predefined responses (Giovis & Rozsa, 2023). AI-driven processes like chatbots and robotic process automation (RPA) offer opportunities to enhance service desk levels. Self-service logging via online portals and web apps can automate common inquiries, reducing telephone contact and freeing staff to handle complex issues, leading to improved service levels with reduced wait times and 24-hour availability (Steinig, 2023).

Generative AI, leveraging large language models and deep learning techniques, has the potential to revolutionize customer service by understanding intricate inquiries and generating natural conversation responses (Giovis & Rozsa, 2023). These models improve customer service by analyzing conversations, generating contextually relevant responses, and managing complex queries. They understand nuanced intent, sentiment, and context, providing personalized answers, recommendations, and solutions, thus enhancing the overall customer experience. IBM identifies five key use cases where generative AI can disrupt the customer service industry (Giovis & Rozsa, 2023):

1. **Conversational Search:** Generative AI uses finely tuned language models with company knowledge bases to provide natural, relevant responses in the user's preferred language, reducing the need for translation services and minimizing search effort.
2. **Agent Assistance – Search & Summarization:** Generative AI helps customer support agents respond to inquiries by generating automatic responses in the preferred communication channel and creating concise summaries that employees can use to offer product information, services, or recommendations, while categorizing and tracking trends.
3. **Build Assistance:** Generative AI aids in content creation and building service tools to process requests, generating responses and suggestions based on company or customer data.
4. **Call Center Operational Data Optimization:** Generative AI can handle repetitive tasks to enhance the feedback loop within a call center.
5. **Personalized Recommendations:** By considering a customer's interaction history with a brand, generative AI provides information specifically relevant to that customer.

Overall, the integration of generative AI in customer service can significantly enhance the efficiency and quality of interactions, offering more natural and personalized responses and improving the overall customer experience.

10 Acknowledgments

The creation of this paper was supported by Benjamin

Ledel, the owner of Digital Learning GmbH and LLM educa AI, who volunteered as interview partner.

11 Conclusion

In conclusion, the deployment and development of LLMs illustrate a significant leap in NLU and NLP. Technological innovations such as GenRec LLM, leveraging raw text data for personalized recommendations, and LoRa's low-rank matrix approximations highlight the evolving landscape of generative recommendation systems (Cabello, 2023; Ji et al., 2024). The rise of smaller, more efficient models like LLaMA and the proliferation of open-source LLMs such as XLNET and Qwen 1.5 democratize access to advanced AI capabilities, ensuring broader dissemination and application (Huynh, 2024; Touvron, Lavril, et al., 2023; Touvron, Martin, et al., 2023).

The contributions of pioneering researchers and organizations have been instrumental in advancing LLM technology. Visionaries like Turing, Weizenbaum, and Winograd laid foundational principles, while contemporary figures like Hinton and Karpathy have refined training methods and architectures (Cabello, 2023; Gonçalves, 2023; Rusch, 2023; Tarnoff, 2023). Companies like Google, Meta, OpenAI, and innovative startups like Hugging Face continue to push the boundaries of what LLMs can achieve.

Despite their potential, LLMs face significant challenges, including the ethical implications of misinformation, plagiarism, and the societal impacts on the job market (Agarwal et al., 2023; Arntz et al., 2016; Mündler et al., 2024). Addressing these issues requires robust policies, ethical guidelines, and international cooperation to ensure responsible usage.

In practical applications, LLMs are transforming fields such as education, law, and customer service. Tools like Google Translate and Grammarly, and initiatives like educa AI and Syntea, exemplify how LLMs can enhance learning and professional environments (Gimpel et al., 2023). In Swiss German contexts, projects like ASR by FHNW and Swisscom's Whisper demonstrate efforts to overcome linguistic challenges (Radford et al., 2023; University of Applied Sciences and Arts of Northwestern Switzerland, n.d.).

Overall, the advancement of LLMs in the DACH region presents both opportunities and challenges, necessitating a balanced approach to harness their benefits while mitigating potential risks. As research and development continue, the potential for LLMs to revolutionize various sectors remains vast and promising.

12 Outlook

The democratization of LLMs by startups like Hugging Face has made these tools more accessible, emphasizing the importance of collaboration among researchers, developers, and institutions to accelerate progress and enhance their effectiveness and reach.

As LLMs become more prevalent across various industries, addressing ethical and societal issues becomes paramount. It is crucial to ensure fairness, transparency in response generation, and to tackle privacy concerns.

Responsible development and deployment practices are essential to mitigate these challenges, ensuring that LLMs have a positive impact on society.

In the educational sector, the importance of context-based answers cannot be overstated, as generic responses may fail to reflect the specific material presented by educators. Moreover, the development of LLMs must overcome the unique challenges posed by languages such as German and Swiss German. Addressing these linguistic challenges is essential for the effective application of LLMs within the Swiss German context.

The Swiss market presents a high potential for the adoption of advanced AI tools, particularly in corporate and educational sectors. Companies offering these tools should consider collaborating with Swiss universities, taking advantage of the current absence of significant competitors in the market. Such collaborations can drive innovation and ensure that LLMs are tailored to meet the specific needs of the Swiss market, both in terms of language and application.

By focusing on these areas, we can optimize the development and deployment of LLMs, ensuring they are ethical, effective, and widely accessible, ultimately driving significant advancements in AI and NLP.

In the next phase of this research, we will engage in fine-tuning an LLM using self-gathered data input. This will involve training the model with data specific to our research focus, allowing for a more tailored application of the LLM. The findings from this fine-tuning process will be detailed in the subsequent part of this paper, providing insights into the model's performance, its ability to handle context-specific queries, and its efficacy in addressing the linguistic challenges particular to Swiss German. This hands-on approach will contribute to the broader understanding of LLM customization and its practical applications in various contexts.

By focusing on these areas, we can optimize the development and deployment of LLMs, ensuring they are ethical, effective, and widely accessible, ultimately driving significant advancements in AI and natural language processing.

REFERENCES

- Abadi, M., Yu, Y., Barham, P., Brevdo, E., Burrows, M., Davis, A., Dean, J., Ghemawat, S., Harley, T., Hawkins, P., Isard, M., Kudlur, M., Monga, R., Murray, D., & Zheng, X. (2018). Dynamic Control Flow in Large-Scale Machine Learning. *Proceedings of the Thirteenth EuroSys Conference*, 1–15. <https://doi.org/10.1145/3190508.3190551>
- Agarwal, A., Padhi, A., Vala, J., & Katoch, C. D. S. (2023). Large Language Models in Academia: Ethical Considerations and Future Prospects. *NMO Journal*, 17(2), 105–106. https://doi.org/10.4103/JNMO.JNMO_16_23
- Alomari, A., Idris, N., Sabri, A. Q. M., & Alsmadi, I. (2022). Deep reinforcement and transfer learning for abstractive text summarization: A review.

- Computer Speech & Language*, 71, 101276. <https://doi.org/10.1016/j.csl.2021.101276>
- Arntz, M., Gregory, T., & Zierahn, U. (2016). *The Risk of Automation for Jobs in OECD Countries: A Comparative Analysis* (189). OECD. <https://doi.org/10.1787/5jlz9h56dvq7-en>
- Bahdanau, D., Cho, K., & Bengio, Y. (2016). *Neural Machine Translation by Jointly Learning to Align and Translate* (arXiv:1409.0473). arXiv. <https://doi.org/10.48550/arXiv.1409.0473>
- BaiDoo-Anu, D., & Owusu Ansah, L. (2023). Education in the Era of Generative Artificial Intelligence (AI): Understanding the Potential Benefits of ChatGPT in Promoting Teaching and Learning. *Journal of AI*, 7(1), 52–62. <https://doi.org/10.61969/jai.1337500>
- Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? 🦜. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623. <https://doi.org/10.1145/3442188.3445922>
- Berner Fachhochschule. (n.d.-a). *Legal NLP - Der Einsatz von grossen und kleinen Sprachmodellen im Rechtswesen*. Retrieved December 5, 2024, from [https://www.bfh.ch/dam/jcr:a3b4b524-74d8-4fa6-b287-e8f2abaa38c5/2024_Legal%20NLP%20\(Public%20Sector%20Perspectives\).pdf](https://www.bfh.ch/dam/jcr:a3b4b524-74d8-4fa6-b287-e8f2abaa38c5/2024_Legal%20NLP%20(Public%20Sector%20Perspectives).pdf)
- Berner Fachhochschule. (n.d.-b). *Legal NLP: Datensätze und Publikationen zu Künstlicher Intelligenz im Rechtswesen*. Retrieved May 12, 2024, from <https://www.bfh.ch/de/forschung/forschungsbereiche/public-sector-transformation/kompetenzen/kuenstliche-intelligenz-und-grosse-sprachmodelle-im-rechtswesen-ressourcen-und-publikationen/>
- Bradshaw, S., & Howard, P. N. (2019). *The Global Disinformation Order: 2019 Global Inventory of Organised Social Media Manipulation*. Project on Computational Propaganda.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., ... Amodei, D. (2020). *Language Models are Few-Shot Learners* (arXiv:2005.14165). arXiv. <http://arxiv.org/abs/2005.14165>
- Burrell, J. (2016). How the machine ‘thinks’: Understanding opacity in machine learning algorithms. *Big Data & Society*, 3(1), 205395171562251. <https://doi.org/10.1177/2053951715622512>
- Cabello, A. (2023, November 2). *The Evolution of Language Models: A Journey Through Time*. <https://medium.com/@adria.cabello/the-evolution-of-language-models-a-journey-through-time-3179f72ae7eb>
- Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., Chen, H., Yi, X., Wang, C., Wang, Y., Ye, W., Zhang, Y., Chang, Y., Yu, P. S., Yang, Q., & Xie, X. (2024). A Survey on Evaluation of Large Language Models. *ACM Transactions on Intelligent Systems and Technology*, 15(3), 1–45. <https://doi.org/10.1145/3641289>
- Deci Research team. (2024, March 27). Top Large Language Models Reshaping the Open-Source Arena. *DECI*. <https://deci.ai/blog/list-of-large-language-models-in-open-source/#Qwen1.5-anchor>
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding* (arXiv:1810.04805). arXiv. <http://arxiv.org/abs/1810.04805>
- educa AI. (2021, April 23). *Übersicht—Educa AI*. <https://www.educaai.de/>
- ETH AI Center & EPFL AI Center. (n.d.). *Leveraging the world's most AI-capable supercomputer*. Swiss AI. Retrieved May 11, 2024, from <https://www.swiss-ai.org>
- European Commission. (2022). *Ethical guidelines on the use of artificial intelligence (AI) and data in teaching and learning for educators*. Publications Office of the European Union. <https://data.europa.eu/doi/10.2766/153756>
- Gan, W., Qi, Z., Wu, J., & Lin, J. C.-W. (2023). *Large Language Models in Education: Vision and Opportunities* (arXiv:2311.13160). arXiv. <http://arxiv.org/abs/2311.13160>
- Gerlach, J., Bouillon, P., Vázquez, S. R., Mutal, J., & Starlander, M. (2023). Evaluating a Multilingual Pre-trained Model for the Automatic Standard German captioning of Swiss German TV. In H. Ghorbel, M. Sokhn, M. Cieliebak, M. Hürlimann, E. de Salis, & J. Guerne (Eds.), *Proceedings of the 8th edition of the Swiss Text Analytics Conference* (pp. 14–22). Association for Computational Linguistics. <https://aclanthology.org/2023.swisstext-1.2>
- Gimpel, H., Hall, K., Decker, S., Eymann, T., Lämmermann, L., Mädche, A., Röglinger, M., Ruiner, C., Schoch, M., Schoop, M., Urbach, N., & Vandirck, S. (2023). *Unlocking the Power of Generative AI Models and Systems such as GPT-4 and ChatGPT for Higher Education A Guide for Students and Lecturers*. <https://doi.org/10.13140/RG.2.2.20710.09287/2>

- Giovis, R., & Rozsa, E. (2023, July 17). *Transforming customer service: How generative AI is changing the game—IBM Blog*. <https://www.ibm.com/blog/transforming-customer-service-how-generative-ai-is-changing-the-game/>
- Gonçalves, B. (2023). Can machines think? The controversy that led to the Turing test. *AI & SOCIETY*, 38(6), 2503–2509. <https://doi.org/10.1007/s00146-021-01318-6>
- Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- Hochschule Luzern – Design Film Kunst. (n.d.). *Whisper OpenAI – Transkription*. Retrieved May 11, 2024, from <https://sites.hslu.ch/werkstatt/whisper-openai-transkription/>
- Honnet, P.-E., Popescu-Belis, A., Musat, C., & Baeriswyl, M. (2017). *Machine Translation of Low-Resource Spoken Dialects: Strategies for Normalizing Swiss German*. <https://doi.org/10.48550/ARXIV.1710.11035>
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., & Chen, W. (2021). *LoRA: Low-Rank Adaptation of Large Language Models*. <https://doi.org/10.48550/ARXIV.2106.09685>
- Huynh, D. (2024, January 27). State Of LLM In 2023: A Quick Recap On Latest Advancements. *Medium*. <https://medium.com/@vndee.huynh/state-of-llm-in-2023-a-quick-recap-on-latest-advancements-46a55dfe1fe5>
- IBM. (n.d.). *What is LLAMA?* IBM. Retrieved May 11, 2024, from <https://www.ibm.com/topics/llama-2>
- IU Internationale Hochschule. (n.d.-a). *Persönlicher Lernassistent Syntea mit KI | IU Internationale Hochschule*. IU – Internationale Hochschule. Retrieved May 4, 2024, from <https://www.iu.de/syntea/>
- IU Internationale Hochschule. (n.d.-b). *Syntea: KI-Assistent verändert die Online-Bildung | IU News*. Syntea: KI-Assistent verändert die Online-Bildung | IU News. Retrieved May 4, 2024, from <https://www.iu.de/news/syntea-ki-gestuetzte-loesung-revolutioniert-online-bildung-und-interaktion-mit-studierenden/>
- Javaid, M., Haleem, A., Singh, R. P., Khan, S., & Khan, I. H. (2023). Unlocking the opportunities through ChatGPT Tool towards ameliorating the education system. *BenchCouncil Transactions on Benchmarks, Standards and Evaluations*, 3(2), 100115. <https://doi.org/10.1016/j.tbench.2023.100115>
- Ji, J., Li, Z., Xu, S., Hua, W., Ge, Y., Tan, J., & Zhang, Y. (2024). GenRec: Large Language Model for Generative Recommendation. In N. Goharian, N. Tonello, & Y. He, A. Lipani, G. McDonald, C. Macdonald, & I. Ounis (Eds.), *Advances in Information Retrieval* (Vol. 14610, pp. 494–502). Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-56063-7_42
- Kasneci, E., Sessler, K., Küchemann, S., Bannert, M., Dementieva, D., Fischer, F., Gasser, U., Groh, G., Günnemann, S., Hüllermeier, E., Krusche, S., Kutyniok, G., Michaeli, T., Nerdel, C., Pfeffer, J., Poquet, O., Sailer, M., Schmidt, A., Seidel, T., ... Kasneci, G. (2023). ChatGPT for good? On opportunities and challenges of large language models for education. *Learning and Individual Differences*, 103, 102274. <https://doi.org/10.1016/j.lindif.2023.102274>
- Khurana, D., Koli, A., Khatter, K., & Singh, S. (2023). Natural language processing: State of the art, current trends and challenges. *Multimedia Tools and Applications*, 82(3), 3713–3744. <https://doi.org/10.1007/s11042-022-13428-4>
- Köchli, O., Wenk, P., Zweili, C., & Hanne, T. (2023). Using BERT for Swiss German Sentence Prediction. In J. Abawajy, J. Tavares, L. Kharb, D. Chahal, & A. B. Nassif (Eds.), *Information, Communication and Computing Technology* (Vol. 1841, pp. 3–15). Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-43838-7_1
- Kohne, A., Kleinmanns, P., Rolf, C., & Beck, M. (2020). *Chatbots: Aufbau und Anwendungsmöglichkeiten von autonomen Sprachassistenten*. Springer Fachmedien Wiesbaden. <https://doi.org/10.1007/978-3-658-28849-5>
- Kraus, S., Kanbach, D. K., Krysta, P. M., Steinhoff, M. M., & Tomini, N. (2022). Facebook and the creation of the metaverse: Radical business model innovation or incremental transformation? *International Journal of Entrepreneurial Behavior & Research*, 28(9), 52–77. <https://doi.org/10.1108/IJEBR-12-2021-0984>
- Kreps, S. (2020). *Replication Data for: All the News that's Fit to Fabricate: AI-Generated Text as a Tool of Media Misinformation* [dataset]. Harvard Dataverse. <https://doi.org/10.7910/DVN/1XVYU3>
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444. <https://doi.org/10.1038/nature14539>
- Microsoft. (2024, April 23). *Introducing Phi-3: Redefining what's possible with SLMs*. <https://azure.microsoft.com/en-us/blog/introducing-phi-3-redefining-whats-possible-with-slms/>
- Mollaki, V. (2024). Death of a reviewer or death of peer review integrity? The challenges of using AI tools in peer reviewing and the need to go beyond publishing policies. *Research Ethics*, 20(2), 239–

250. <https://doi.org/10.1177/17470161231224552>
- Mukherjee, S., Mitra, A., Jawahar, G., Agarwal, S., Palangi, H., & Awadallah, A. (2023). *Orca: Progressive Learning from Complex Explanation Traces of GPT-4* (arXiv:2306.02707). arXiv. <http://arxiv.org/abs/2306.02707>
- Mündler, N., He, J., Jenko, S., & Vechev, M. (2024). *Self-contradictory Hallucinations of Large Language Models: Evaluation, Detection and Mitigation* (arXiv:2305.15852). arXiv. <https://doi.org/10.48550/arXiv.2305.15852>
- Niklaus, J., Chalkidis, I., & Stürmer, M. (2021). *Swiss-Judgment-Prediction: A Multilingual Legal Judgment Prediction Benchmark*. <https://doi.org/10.48550/ARXIV.2110.00806>
- Paonessa, C., Schraner, Y., Deriu, J., Hürlimann, M., Vogel, M., & Cieliebak, M. (2023). *Dialect Transfer for Swiss German Speech Translation* (arXiv:2310.09088). arXiv. <http://arxiv.org/abs/2310.09088>
- Plüss, M., Neukom, L., Scheller, C., & Vogel, M. (2020). *Swiss Parliaments Corpus, an Automatically Aligned Swiss German Speech to Standard German Text Corpus*. <https://doi.org/10.48550/ARXIV.2010.02810>
- Radford, A., Kim, J. W., Brockman, G., McLeavey, C., & Sutskever, I. (2023). *Robust Speech Recognition via Large-Scale Weak Supervision*. OpenAI. <https://cdn.openai.com/papers/whisper.pdf>
- Rizvi, M. S. Z. (2019, August 10). A Comprehensive Guide to Build your own Language Model in Python! *Analytics Vidhya*. <https://medium.com/analytics-vidhya/a-comprehensive-guide-to-build-your-own-language-model-in-python-5141b3917d6d>
- Rusch, L. (2023, January 6). Der „Ketzer der Informatik“: Was ChatGPT und Internet mit Joseph Weizenbaum zu tun haben. *Der Tagesspiegel Online*. <https://www.tagesspiegel.de/wissen/der-ketzer-der-informatik-was-chatgpt-und-internet-mit-joseph-weizenbaum-zu-tun-haben-9137928.html>
- Shokri, R., Song, L., & Mittal, P. (2019). Privacy Risks of Securing Machine Learning Models against Adversarial Examples. *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 241–257. <https://doi.org/10.1145/3319535.3354211>
- Slats, L. (n.d.). *A Brief History of LoRa: Three Inventors Share Their Story*. Retrieved May 25, 2024, from <https://blog.semtech.com/a-brief-history-of-lora-three-inventors-share-their-personal-story-at-the-things-conference>
- Spitch. (n.d.). *Spitch*. Start Now. Retrieved May 11, 2024, from <https://spitch.ai/de/>
- Spitch. (2020, July 13). *Migros Bank wählt Spitch-Lösung zur Identifizierung ihrer Kunden*. <https://spitch.ai/de/news/2020.07.13.page>
- Steinig, T. (2023). *Meatred OpenAI & LLM- Disruption im Bereich des Service Desk* [Projektarbeit, Hochschule St. Gallen]. <https://sgbs.ch/wp-content/uploads/Projektarbeit-TSteinig.pdf>
- SwissLLM. (n.d.). SwissLLM. Retrieved May 11, 2024, from <https://swissllm.ch/>
- Tarnoff, B. (2023, July 25). Weizenbaum's nightmares: How the inventor of the first chatbot turned against AI. *The Guardian*. <https://www.theguardian.com/technology/2023/jul/25/joseph-weizenbaum-inventor-eliza-chatbot-turned-against-artificial-intelligence-ai>
- Thommen, K. (2023). *Swiss german Speech-to-Text: Test and Improve the Performance of Model on Spontaneous Speech* [Master Thesis, University of Zurich]. <https://www.merlin.uzh.ch/contributionDocument/download/16268>
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., & Lample, G. (2023). *LLaMA: Open and Efficient Foundation Language Models*. <https://doi.org/10.48550/ARXIV.2302.13971>
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C. C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., ... Scialom, T. (2023). *Llama 2: Open Foundation and Fine-Tuned Chat Models*. <https://doi.org/10.48550/ARXIV.2307.09288>
- Universität Zürich. (n.d.). *Teaching Tools*. Artificial Intelligence. Retrieved May 4, 2024, from <https://teachingtools.uzh.ch/en/tools/kuenstliche-intelligenz>
- University of Applied Sciences and Arts of Northwestern Switzerland. (n.d.). *Speech Recognition for Swiss German*. FHNW. Retrieved May 11, 2024, from <https://www.fhnw.ch/en/about-fhnw/schools/school-of-engineering/institutes/research-projects/speech-recognition-for-swiss-german>
- Vamvas, J., Graen, J., & Sennrich, R. (2023). *SwissBERT: The Multilingual Language Model for Switzerland*. <https://doi.org/10.48550/ARXIV.2303.13310>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2023). *Attention Is All You Need* (arXiv:1706.03762). arXiv. <http://arxiv.org/abs/1706.03762>
- Walther, M. (2024, January 23). *Students use language models primarily for programming*. Students Use Language Models Primarily for Programming. <https://ethz.ch/staffnet/en/news-and->

events/internal-news/archive/2024/01/students-use-language-models-primarily-for-programming.html

Widmer, C. (2018, June 4). *Algorithmus lernt Schweizerdeutsch* / *Swisscom*. <https://www.swisscom.ch/de/business/enterprise/themen/digital-business/spracherkennung-schweizerdeutsch.html>

Xayn. (n.d.). *Noxtua Copilot*. Noxtua. Retrieved May 12, 2024, from <https://www.noxtua.ai>

Yan, L., Sha, L., Zhao, L., Li, Y., Martinez-Maldonado, R., Chen, G., Li, X., Jin, Y., & Gašević, D. (2024). Practical and ethical challenges of large language models in education: A systematic scoping review. *British Journal of Educational Technology*, 55(1), 90–112. <https://doi.org/10.1111/bjet.13370>

Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., Du, Y., Yang, C., Chen, Y., Chen, Z., Jiang, J., Ren, R., Li, Y., Tang, X., Liu, Z., ... Wen, J.-R. (2023). *A Survey of Large Language Models* (arXiv:2303.18223). arXiv. <http://arxiv.org/abs/2303.18223>

Seraina Fischer BSc Business IT, 2021. Migros Bank.

Bigna Schmid BSc International Management, 2021. UBS Business Solutions AG.

Annamària Sréter BSc Agricultural Sciences, 2019. Resilux Schweiz AG.