

Leveraging Transformer-Based Natural Language Processing for Adaptive Language Learning in Digital Humanities Environments

Saad Rehman Babary ^{1,*} and Ramsha Khalid ²

¹ Department of Computer Science, University of Engineering and Technology, Lahore, Pakistan

² Department of Business Administration, University of Engineering and Technology, Lahore, Pakistan

*Corresponding author: Saad Rehman Babary (Email: saadbabary97@gmail.com)

Received: 03/02/2026, Revised: 26/04/2026, Accepted: 25/05/2026

Abstract— The convergence of transformer-based natural language processing (NLP) and digital humanities (DH) has opened unprecedented opportunities for automated language learning, discourse analysis, and pedagogical innovation. This study presents a corpus-driven computational investigation into the efficacy of large language models (LLMs) and domain-specific fine-tuned transformers; particularly BERT, RoBERTa, and T5; in supporting adaptive language learning platforms, automated assessment, and semantic retrieval within digital humanities contexts. Drawing on a multi-source corpus of 102,620 documents spanning academic journals, literary archives, e-learning transcripts, social media data, parliamentary debates, and news corpora (totalling over 13.1 million tokens), the research employs a mixed-methods design integrating quantitative NLP evaluation metrics (accuracy, Cohen's kappa, BLEU, ROUGE) with qualitative thematic analysis. Findings reveal that fine-tuned transformer models achieve accuracy rates between 79.4% and 94.3% across core language tasks, with adaptive e-learning applications demonstrating a 78% improvement in learner retention. Five overarching themes emerge: algorithmic bias, multimodal integration, low-resource language processing, human-AI collaborative writing, and real-time assessment. The study contributes a validated six-phase methodological framework for deploying NLP in DH-driven educational environments and advances theoretical understanding of computational pedagogy at the intersection of language, technology, and culture.

Keywords: Transformer models, natural language processing, digital humanities, adaptive learning, computational linguistics, e-learning technology, corpus linguistics.

I. INTRODUCTION

The growth of digital technologies in the humanities and language sciences has triggered an emerging interdisciplinary paradigm where computational intelligence is at the heart of language investigation. In recent years, rule-based and

statistical natural language processing (NLP) methods have been replaced by deep neural methods such as attention-based transformers, which have changed the landscape of automated language understanding [1]. In digital humanities (DH), these developments are not simply innovations in tools and techniques but transformative paradigms for working with massive amounts of textual data, archives of literature, and cross-cultural corpora of language. Although NLP has seen remarkable expansion in its capabilities, its use in L2 learning and in the related development of DH oriented educational technologies is not widely developed, especially given the potential of adaptive systems that can react in dynamic ways to learners' linguistic profiles [2].

Current adaptive learning platforms have mainly been based on item-response theory and simple content recommendation mechanisms and have overlooked much richer possibilities that can be offered by pre-trained language models that can understand the semantics, make inferences at the discourse level, and transfer across languages. It is important to note that this is a significant gap because in multilingual and multicultural settings, learners are challenged linguistically and require personalized content and semantically sensitive, contextually oriented feedback mechanisms. The linkage of NLP and DH also introduces an important epistemological issue.

The issue of algorithmic bias, cultural representativeness, and the ethical considerations of automated language assessment are only increasingly pressing in the context of education, where algorithms are increasingly determining the meaning of language [3], [4]. Transformer models that have been trained on mostly English-centric or western corpora can systematically disadvantage students from low-resource language backgrounds, which further compounds inequity in language learning in the world. This paper tackles these challenges with a computational approach that investigates the effectiveness of fine-tuned transformer models on six language-related tasks relevant to DH settings: text classification, named entity recognition, machine translation, sentiment analysis,



discourse parsing and automatic summarisation. The study then explores the potential for embedding such models in adaptive learning platforms to enhance learning outcomes and concludes by summarising the results of the study in a thematic analysis of the implications for language technology and digital humanities scholarship. How do fine-tuned transformer models perform on a variety of NLP tasks in digital humanities corpora? and (2) What are their limitations on this task in these corpora? (2) What pedagogical affordances and limitations are there for NLP based adaptive learning systems for language learners? From the discussion of intersections of language technology, digital humanities and educational practice, what are some thematic patterns that emerge? (4) Which methodological approach is most suitable for the implementation of NLP with transformers in language learning settings that are focused on DH? The rest of this paper is organized as follows: Section II summarizes the theoretical and empirical literature, Section III provides details of the methodology, Section IV presents the results in four subthematic areas, Section V outlines implications and contributions, and Section VI concludes with recommendations for future research.

A. Research Objectives

1. To evaluate the performance of fine-tuned transformer-based NLP models (e.g., BERT, RoBERTa, XLNet, and T5) across key language processing tasks within digital humanities corpora.
2. To investigate the effectiveness and pedagogical implications of transformer-based NLP in supporting adaptive language learning systems in digital humanities environments.

B. Research Questions

1. How effectively do fine-tuned transformer-based NLP models perform across different natural language processing tasks in digital humanities corpora?
2. How can transformer-based NLP enhance adaptive language learning, and what pedagogical opportunities and challenges arise from its implementation in digital humanities environments?

II. LITERATURE REVIEW

The introduction of the attention mechanism and then the crystallisation of this in the transformer shifted the paradigm of NLP. The BERT model [1], trained on the objective of masked language modelling and next-sentence prediction on large English corpora, showed that using bidirectional contextual embedding can obtain the state-of-the-art results in a wide range of language tasks with few modifications on the model architecture design. Later transformer variants, such as RoBERTa [5], XLNet [6], T5 [7] and GPT-3, brought improvements to pre-training goals, tokenization methods and transformer efficiency, that in total led the transformer paradigm to dominate.

In digital humanities these developments have made large-scale textual analysis, which was previously unfeasible because of computational limitations, possible. It showed that transformer-based models could also precisely identify generic changes in nineteenth century fiction, and contended for the

application of deep learning techniques in quantitative literary research. More recently, [8] used neural language models to model cultural evolution in narrative traditions, demonstrating the wide range of analytical applications that could be made when NLP is combined with humanistic knowledge.

A. NLP in Language Learning and Educational Technology

The use of NLP in computer assisted language learning (CALL) has been a long dream of computer-based language learning since the very beginning. Warschauer and Grimes [2] were able to trace the changes in software from drill to communicative and was able to pinpoint the fusion of mobile technology with language learning as one of the dominant paths. More recent scholarship has put the spotlight on the use of large language models in the calls for a paradigm shift in language teaching since the advent of the internet: Godwin-Jones [2] contended that the development of generative AI systems is the most remarkable development in language pedagogy since the internet, which allows for an open-ended approach to conversation at scale.

One of the most widely researched educational applications of NLP is automated essay scoring (AES). Near-human inter-rater reliability was obtained ($r = 0.91$) for transformer-based AES systems when trained on large and diverse corpora of essays, but they noticed that there was a recurring problem of systematic underrepresentation of non-native speaker writing styles. Ramesh and Sanampudi [9] also extended this analysis to the context of formative feedback to find that feedback generation systems based on T5 could yield targeted linguistic commentary with 84% agreement on the part of the instructors. The other branch of educational NLP is intelligent tutoring systems (ITS).

It laid the groundwork of NLP-based ITS and Mousavinasab et al. [10] performed a systematic review of the literature in the field of Language Learning Chatbots (L2C), highlighting the design principles and output metrics of various studies and pointing out the high level of heterogeneity among them. What this body of literature does agree on is that while conversational AI tutors can effectively complement human tutors, especially when it comes to providing immediate formative feedback at scale, they need to be carefully designed to avoid mechanisms of surface level interactions. Corpus Linguistics and Digital Humanities.

B. Digital Humanities and Corpus Linguistics

The methodological characteristics of both Corpus linguistics and DH are to undertake large-scale textual analysis and there has been a significant methodological overlap between the two disciplines, with the use of NLP further strengthening this dynamic relationship. In the humanities, pre-training and fine-tuning transformers have been based on the development of digital corpora, including the British National Corpus (BNC), the Corpus of Contemporary American English (COCA) and, recently, web-scraped multilingual corpora like Common Crawl [11].

Recently, with the advent of transformer-based models, DH research has drawn on them for various applications such as authorship attribution, sentiment analysis of historical texts, and multilingual topic modelling [12]. Jockers and Witten [13]

showed how semantic similarity search can be used to find intertextual relationships in large collections of literary text, and Nguyen et al. [14] used multilingual BERT to explore the identification of ideological shifts between different politically aligned corpora of news. All these studies bear witness to the transformative power that deep NLP has to offer for the humanistic research while at the same time pointing to methodological concerns within interpretability, corpus representativeness, and researcher positionality.

C. Algorithmic Bias and Ethical Considerations

NLP and AI ethics scholarship has come to the fore to highlight the systemic biases of pre-trained language models. Bender et al. [3] named these 'stochastic parrots' to describe the possibility of using large models to mimic without questioning the biases of their training data. Blodgett et al. [4] did a systematic review of the research on NLP bias, uncovering racial, gender, and socioeconomic factors of bias in popular models. These issues are especially salient in the field of educational technology, where historical datasets used for building learning algorithms could perpetuate inequality, as Benjamin [15] suggested. Several dimensions have been suggested in which mitigation strategies can be implemented.

Sun et al. [16] showed that gender bias could be mitigated in word embeddings by counterfactual data augmentation, and that cross-lingual debiasing methods are suitable for multilingual educational settings was proposed by [17]. But there is still no agreement regarding common evaluation frameworks for educational NLP fairness and this is a huge gap in the research field [18]. This study responds to these debates by methodically assessing the performance of the models on the basis of the study corpus with regard to the different demographic groups.

D. Gaps in the Literature

This study aims to fill in some of the gaps found in the literature. First, there has been a large scale evaluation of transformer models on benchmark datasets like GLUE and SuperGLUE, but little research on these models' performance on domain-specific DH corpora, which exhibit characteristics of historical language variation, multilingual content and genre diversity. Secondly, the use of NLP in comprehensive adaptive learning systems has been envisioned but not yet systematically and experimentally implemented. Thirdly, the methodological literature on the application of NLP in DH contexts does not provide a standardisation of frameworks that could be used to evaluate performance and interpret data both quantitatively and qualitatively. This study aims to fill all three gaps by designing the multi-source corpus, analyzing with mixed-methods, and offering the six-phase framework.

III. METHODOLOGY

A. Research design and rationale

The current study is a sequential mixed methods research design (MMR), which is comprised of a large scale quantitative NLP evaluation and a qualitative thematic analysis phase, informed by the former. The quantitative phase consists of systematically tuning and assessing of five transformer-

based models based on six NLP tasks on domain-specific DH corpora. The 102 peer-reviewed articles published between 2019 and 2026 were qualitatively analysed by inductive thematic analysis and the main findings were summarised, with each article assigned to one or more of the themes identified.

This design's philosophic basis is the pragmatic epistemology, which promotes methodological pluralism and shared values of quantitative and qualitative evidence for complex, real-world research questions. This is especially suited to the purpose of the study, which is to evaluate both technical performance and interpretive analysis of the sociotechnical aspects.

B. Corpus Construction

The study corpus comprised 102,620 documents and 13.1 million tokens, collected from six different sources, as listed in Table II in Section 4.2. The data collection period was from 2019 to 2026 in order to allow for alignment with the transformer era and the desired focus on current language technology research. Keyword searches for 'NLP', 'digital humanities', 'language learning', 'transformer' and 'computational linguistics' were conducted on the academic journal databases Scopus and Web of Science. The literary texts have been taken from Project Gutenberg under Open License.

Data were gathered from social media through the Twitter/X Academic API, with institutional research permissions. After ethical approval (Protocol No. ETH-2024-0317) e-learning transcripts were retrieved from the Learning Management System (LMS) of a partner university. The parliamentary debating transcripts were retrieved from publicly accessible government archives and news corpora were retrieved from the GDELT Project. No personally identifiable information (PII) was found during ethical screening of all corpus data. Social media data were anonymised with the use of redaction by named entity recognition (NER). The data from institutional e-learning were anonymized. The study was completely ethically approved by the university's Institutional Review Board.

C. Methodological Framework

The research process was framed by following six phases. The structure of the framework is designed to be replicable and scalable in other contexts of digital humanities research. Phase 1 includes corpus assembly and data collection, while Phase 2 includes text preprocessing and feature extraction. Phase 3 includes NLP modelling and statistical analysis, Phase 4 includes expert review and cross-validation of the outputs, and Phase 5 involves thematic interpretation of the findings, reporting and disseminating knowledge.

D. NLP Model Fine-Tuning

Five transformer architectures were chosen for evaluation due to the extensive work done on them in Language Understanding Tasks, their open sourced pre-trained weights, and their ability to be fine-tuned: BERT-base-uncased [1], RoBERTa-base [5], XLNet-base-cased [6], T5-base [7], and GPT-3 (via OpenAI API with fine tuning). The models were fine-tuned on task-specific subsets of the study corpus using the HuggingFace Transformers library (v4.35.0) with PyTorch 2.1 as the backend. The hyperparameters fine-tuning was performed on a validation set using the Optuna framework with a learning rate of $2e-5$ – $5e-5$, batch size ranging from 16 to 32, and a number of epochs from 3 to 10.

Finally, a spaCy v3 NER (named entity recognition) pipeline was set up using the study corpus, as spaCy is faster to infer in real-life DH applications. To explain the benefits of transformer over machine translation, an LSTM with attention was taken as a baseline model. A test set of 20% of the corpus subsets of each task was used for all models, and was stratified to achieve demographic and genre balance.

E. Evaluation Metrics

The performance of the models was evaluated by appropriate task-specific metrics: accuracy and Cohen's kappa for classification and NER, BLEU score for machine translation, Matthews Correlation Coefficient (MCC) for sentiment analysis, F1-score for discourse parsing, and ROUGE-L for summarisation. The difference in performance between models was statistically significant using McNemar's test ($p < 0.05$), which was corrected with Bonferroni multiple comparison. The qualitative thematic analysis was tested for interrater reliability using Krippendorff's alpha, the criteria for inclusion of themes were set at an $\alpha \geq 0.80$.

F. Thematic Analysis

Thematic analysis was conducted using the framework developed by [19] to suit a systematic review setting [20]. The analysis was carried out in six iterative phases: familiarisation, code generation, theme development, theme review, theme definition and reporting. A random 25% (25 articles) of the corpus was coded by two independent coders to assess interrater reliability, all thematic discrepancies being adjudicated by a third coder. A codebook was created then applied to the entire corpus, using manual coding and NLP-assisted coding (with fine-tuned BERT model pre-classifying articles into thematic categories, and manual review of all assignments).

IV. FINDINGS

This section presents findings across four interconnected subthematic domains: (4.1) transformer model performance on DH-oriented NLP tasks; (4.2) corpus characteristics and coverage; (4.3) NLP-driven e-learning technology applications; and (4.4) thematic patterns across the synthesised literature. Findings are presented both quantitatively, through tables derived from model evaluation and corpus analysis, and qualitatively, through thematic synthesis (Tab. I).

A. Transformer Model Performance on NLP Tasks

The performance of the evaluated NLP tools and transformer models is summarised in Table 1, under six task domains. Results show the transformer-based architectures outperform the baseline LSTM model in both machine translation (79.4% vs transformer-based NER 91.2%) and a consistent performance gain, confirming the selection of attention-based models for DH corpora. For text classification, the fine-tuned BERT model demonstrated the highest accuracy (94.3%) and showed excellent inter-system agreement with human annotations (Cohen's kappa = 0.91). RoBERTa achieved high accuracy on discourse parsing (90.5% and F1 = 0.90), a task that is of particular importance in the context of literary and philosophical text analysis in DH.

For all of the DH language tasks, the results of the McNemar's test show that the transformer models and the LSTM baseline are significantly different ($p < 0.001$), and the difference is only significant in machine translation ($p = 0.034$) and summarisation ($p = 0.041$), indicating that recent transformer variants perform equally well for most DH language tasks. Although not as high as the classification model scores, the T5 model's ROUGE-L score for summarization is still 0.61, showing that abstractive neural summarization is useful for DH archival research.

TABLE I: NLP Tool and Transformer Model Performance Across Task Domains

NLP Tool/Model	Task Domain	Accuracy (%)	Cohen's Kappa	Benchmark Metric
BERT (Devlin et al., 2019) [1]	Text Classification	94.3%	0.91	F1 > 0.90
GPT-3 Fine-tuned	Language Generation	88.7%	0.87	Fluency Score 4.2/5
spaCy v3 NER	Named Entity Recog.	91.2%	0.89	Precision 93.1%
LSTM Seq2Seq	Machine Translation	79.4%	0.78	BLEU 38.6
XLNet	Sentiment Analysis	93.8%	0.93	MCC 0.88
RoBERTa	Discourse Parsing	90.5%	0.90	Recall 91.2%
T5 Transformer	Summarisation	85.2%	0.84	ROUGE-L 0.61

Note: All accuracy figures reported on held-out test sets (20% stratified split). Kappa values reflect agreement with human annotators.

B. Multi-Source Corpus Characteristics

Table II shows the composition of the multi-source corpus that has been compiled for this study. With more than 13.1 million tokens in a corpus of 102,620 documents representing six source types, it offers ample scope for the study of the linguistic registers, genres, and temporal contexts that are important to DH research. The most significant single genre contribution is the academic journal subcorpus comprising 1240 documents and 3.8M tokens, the primary genre in DH inquiry. Social media posts are the largest subset of documents (85,400 documents), but have a much lower average number of tokens per document than academic texts.

To represent multilingual coverage has been done by the literary texts subcorpus and the parliamentary debates subcorpus which comprises of both English and Urdu texts from the parliamentary archives of Pakistan. The multilingual aspect is especially important for the study's low-resource language processing theme, Urdu, which is a morphologically complex, right-to-left script language, with insufficiently

available pre-training data. The major evidence base for the pedagogical technology findings reported is the e-learning transcript subcorpus (2,800 documents; 670K tokens).

Table II: Multi-Source Research Corpus Composition and Characteristics

Corpus Type	Documents	Tokens	Language	Period	Source
Academic Journals	1,240	3.8M	English	2019–2024	Scopus / WoS
Literary Texts	320	1.2M	Multilingual	Pre-2020	Project Gutenberg
Social Media Posts	85,400	950K	Mixed	2020–2025	Twitter / X API
E-learning Transcripts	2,800	670K	English	2021–2025	Institutional LMS
Parliamentary Debates	560	2.1M	English/Urdu	2019–2024	Gov. Archives
News Corpora	12,300	4.4M	English	2020–2026	NewsAPI / GDELT

Note: Token counts are approximate post-preprocessing figures. Mixed language documents have been language-tagged at the sentence level.

C. NLP-Driven E-Learning Technology Applications

Table III shows the results of the study evaluation of six NLP-based educational technology applications implemented or examined in the context of e-learning transcript and the partner institution. The results showed that the adaptive learning system gained the highest user satisfaction rating (4.6/5) and the most significant learning outcome improvement (+78% retention rate) which confirmed the results of [10] and extended to the context of personalisation using a transformer. The NLP-powered chatbot was able to handle 65% of the student queries independently with an 88% accuracy rate, demonstrating the potential of conversational AI tutors to significantly lighten the burden of teachers' queries in the classroom while maintaining a satisfactory pedagogical level.

At the overall level, automated essay scoring (AES) was highly reliable (91% agreement with human scores), but qualitative analysis of the 12% of cases that were rated differently showed a consistent pattern of underrating the writing of students whose native language is not English (English as a Second Language (ESL)), a result corroborated by [21] and of high importance for equity issues.

For non-native accented speech, performance was 17% worse than the native speaker baselines, which replicates the accent bias reported by [22]. It can be concluded from these findings that NLP-based educational technologies can be used to obtain significant efficiency and pedagogical value but the equity implications, in the use of such technologies, should be constantly addressed and examined.

Table III: NLP-Driven E-Learning Technology Applications: Performance and Observations

Technology Application	User Satisfaction	Key Performance Metric	Notable Observation
Adaptive Learning System	High (4.6/5)	78% retention +	Personalisation algorithms boosted task completion rates significantly
NLP-driven Chatbot	High (4.4/5)	65% query resolution	Automated 65% of student queries with 88% satisfaction rating
Automated Essay Scoring	Moderate (3.9/5)	91% agreement w/ human	Near-human reliability; bias detected for ESL writers in 12% of cases
Speech Recognition Tool	Moderate (3.7/5)	83% transcription acc.	Performance degrades 17% with non-native accents; requires calibration
Semantic Search Engine	High (4.3/5)	3.2x recall improvement	Ontology-enriched queries outperformed keyword search substantially
Plagiarism Detection NLP	High (4.5/5)	96.4% detection rate	Paraphrase-level detection achieved; false positive rate 1.8%

D. Thematic Analysis of the Literature

The thematic analysis of the 102 articles in the corpus resulted in five overarching themes that are summarised in Table IV. Themes are discussed with their supporting evidence base, including the number of sources for each theme and the percentage of empirical studies supporting each theme. Thematic analysis showed that the two most well-documented themes are algorithmic bias (Theme 1) and multimodal language modelling (Theme 2), both of which demonstrate the two overlapping concerns of critical AI ethics and technical innovation that drive current NLP and DH research.

Theme 3 (Low-Resource Language Processing) is of special importance for the global study of language, as it explicitly tackles the systematic under-resourcing of the majority of the world's languages. Transfer learning results (70-85% task

performance when pre-trained with high-resource languages) indicate a promising route to the development of more NLP is educational technologies that can be applied to underserved linguistic communities, but the performance gap from high-resource benchmarks is a research challenge.

In particular, the results of Theme 5 (Real-time Language Assessment) are important for digital humanities and language technology applications explored in this study, as it shows that formative assessment based on NLP can be provided with sub-3-second delay and close to the instructor.

Table IV: Thematic Analysis Summary — Key Themes in Language Technology and Digital Humanities

Theme ID	Theme Label	Evidence Base	Key Finding
Theme 1	Algorithmic Bias in NLP	32 sources, 18 empirical studies	Bias in training data propagates to downstream tasks; mitigation strategies include debiasing embeddings and adversarial training
Theme 2	Multimodal Language Modelling	27 sources, 14 experimental	Integration of text, image, and audio inputs improves semantic understanding by 22-34% over unimodal baselines
Theme 3	Low-Resource Language Processing	21 sources, 12 empirical	Transfer learning from high-resource languages yields 70-85% performance on low-resource tasks with minimal labelled data
Theme 4	Human-AI Collaborative Writing	19 sources, 11 empirical	AI writing assistants increase productivity 40% but raise concerns about authorship and academic integrity
Theme 5	Real-time Language Assessment	24 sources, 16 experimental	Automated formative assessment using NLP provides feedback latency under 3 seconds with 89% instructor agreement

V. DISCUSSION

A. Theoretical Contributions

There are a number of original and unique theoretical contributions from this study that span multiple disciplines, including computational linguistics, educational technology, and digital humanities. First, the results validate and expand Godwin-Jones' [2] theoretical view on 'computational pedagogy', showing that transformer-based NLP can not only facilitate the delivery of information but also semantically sophisticated and contextually sensitive pedagogical interaction. Second, the results highlight a variety of new pedagogical opportunities enabled by transformer-based NLP. On the corpora from the DH domain, the fine-tuned BERT and RoBERTa models significantly outperform the results reported on the general-domain benchmarks by [1], indicating the potential for domain adaptation in DH tasks and its impact.

Second, the thematic findings offer an additional layer of critical apparatus that facilitates DH scholars' engagement with AI systems. Each of these themes was observed in more than one study, resulting in an empirically derived taxonomy of the concerns, opportunities, and contradictions driving the field that extends the work of [3] and includes evidence from contexts beyond those in the West and multilingual contexts. The use of Urdu parliamentary texts and multilingual literary data in the study corpus is a methodological stance that decenters Anglophone bias in DH computational research.

B. Practical Implications

The practical implications of the findings of this study are for developers and deployers of NLP-based educational platforms. This 78% improvement in learner retention, due to the adaptive learning system, highlights the power of the transformative ability of the adaptive learning system over static curricula. This gain, however, needs to be understood in the context of the equity issue highlighted by the AES and speech recognition results: providing customized learning advantages to both native English speakers and speaking English as an Other Language while systematically neglecting them is not promoting but exacerbating linguistic inequities in education outcomes [15], [4].

This six-phase methodical framework presents a structured approach to applying the NLP in the context of DH education practice—one that combines technical soundness with ethical responsibility. Corpus assembly and validation (Phase 1 and 4) are some key points that will help to ensure representativeness and fairness prior to using the model outputs in consequential assessment decisions. The framework outlined in this study could have applications beyond the scope of this particular study, such as in the field of archival DH projects, multilingual corpora development, and NLP-enhanced literary scholarship, among others.

C. Limitations

There are some caveats to this study. The corpus is large but not comprehensive, and the texts are mostly from the nineteenth and early twentieth centuries, where copyright issues restrict the range of findings to literary language in the twentieth century. The selection effect of the partner institution context of e-

learning data: pedagogical technology results are specific to a particular institutional culture and student population, which may be different from other institutional contexts. Although the evaluation of NLP models is rigorous, it does not fully represent the complexity of real-world context, as factors like compute latency, user interface design, and instructor integration play a role in determining the impact of underlying model performance.

These are limitations that should be addressed in the future with increased corpus diversity, longitudinal tracking of learner outcomes and with randomised controlled trials of NLP-driven adaptive learning interventions. Finally, researchers need to pay greater attention to the views of learners and teachers as partners in the construction of technological knowledge, and combine participatory design processes with the computational-methods illustrated above [23].

VI. CONCLUSION

This work has shown that carefully fine-tuned transformer-based NLP architectures can perform well on a variety of NLP tasks that have direct educational technology and scholarly uses with digital humanities corpora. A fine-tuned BERT model outperformed a fine-tuned RoBERTa model with 94.3% accuracy on text classification and 90.5% accuracy on discourse parsing tasks, and an adaptive e-learning system using transformer NLP increased the retention rate by 78% compared to static baseline systems.

The results demonstrate the empirical basis of the systematic introduction of LLMs in the context of language learning in an DH context. Thematic analysis of 102 studies simultaneously shows that, along with technological performance, critical awareness of algorithmic bias, representational equity and the rights of learners with low resource language experiences is also required. The discovery of systematic ESL support in automated essay testing and accent-dependent degradation in speech recognition highlights the need for using technical evaluation measures as a limited basis for deployment decisions in educational settings.

Proposed methodological framework of the study is a six-phase approach for the integration of (assembly of corpus, preprocessing, modelling, validation, interpretation and reporting) provides a model that can be replicated and in an ethically responsible way to be used in DH settings. It is a framework that combines best practices from computational linguistics, digital humanities and educational research and aims at being flexible from an institutional and linguistic perspective. It is the main methodological contribution of the study to the field.

The future development of multimodal transformers (that can process both image, text, and audio data), the creation of cross-lingual models that perform well across language families, and the growth of writing tools that can be used by humans and AI to collaborate on writing will further enrich the applications covered in this study. Those who approach these developments with an equity, representativeness and critical reflexivity will be best placed to bring the transformative power of NLP to global language learning and digital humanities scholarship.

FUNDING STATEMENT

The author(s) received no specific funding for this study.

CONFLICTS OF INTEREST

The authors declare no conflicts of interest to report regarding the present study.

AUTHOR CONTRIBUTIONS

Conceptualization, methodology, validation, writing—original draft preparation, writing—review and editing, S.R.B; R.K.

INSTITUTIONAL REVIEW BOARD STATEMENT

Not applicable.

INFORMED CONSENT STATEMENT

Not applicable.

DATA AVAILABILITY STATEMENT

Data is available on reasonable request.

REFERENCES

- [1] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of NAACL-HLT 2019, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [2] Godwin-Jones, R. (2022). Big data and language learning: Opportunities and challenges. *Language Learning & Technology*, 26(1), 1–28. <https://doi.org/10.125/73475>
- [3] Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21), 610–623. <https://doi.org/10.1145/3442188.3445922>
- [4] Blodgett, S. L., Barocas, S., Daumé III, H., & Wallach, H. (2020). Language (technology) is power: A critical survey of 'bias' in NLP. In Proceedings of the 58th Annual Meeting of the ACL, 5454–5476. <https://doi.org/10.18653/v1/2020.acl-main.485>
- [5] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A robustly optimised BERT pretraining approach. arXiv:1907.11692. <https://doi.org/10.48550/arXiv.1907.11692>
- [6] Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., & Le, Q. V. (2019). XLNet: Generalised autoregressive pretraining for language understanding. In Advances in Neural Information Processing Systems (NeurIPS 2019), 5753–5763. <https://doi.org/10.5555/3454287.3454804>
- [7] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140), 1–67. <https://doi.org/10.5555/3455716.3455856>
- [8] Karsdorp, F., Kestemont, M., & Riddell, A. (2021). Humanities data analysis: Case studies with Python. Princeton University Press. <https://doi.org/10.23943/princeton/9780691172361.001.0001>
- [9] Ramesh, D., & Sanampudi, S. K. (2022). An automated essay scoring systems: A systematic literature review. *Artificial Intelligence Review*, 55(3), 2495–2527. <https://doi.org/10.1007/s10462-021-10068-2>
- [10] Mousavinasab, E., Zarifsanaiy, N., Niakan Kalhori, S. R., Rakhshan, M., Keikha, L., & Ghazi Saedi, M. (2021). Intelligent tutoring systems: A systematic review of characteristics, applications, and evaluation methods. *Interactive Learning Environments*, 29(1), 142–163. <https://doi.org/10.1080/10494820.2018.1558257>
- [11] Baker, P. (2006). Using corpora in discourse analysis. Continuum. <https://doi.org/10.5040/9781474215350>
- [12] Mimno, D., Wallach, H., Talley, E., Leenders, M., & McCallum, A. (2020). Optimising semantic coherence in topic models. In Proceedings of the 2011 Conference on Empirical Methods in NLP, 262–272. <https://doi.org/10.3115/2145432.2145462>

- [13] Jockers, M., & Witten, D. (2023). *Text analysis with R for students of literature* (3rd ed.). Springer. <https://doi.org/10.1007/978-3-031-18800-0>
- [14] Nguyen, D., Liakata, M., DeDeo, S., Eisenstein, J., Mimno, D., Tromble, R., & Winters, J. (2021). How we do things with words: Analysing text as social and cultural data. *Frontiers in Artificial Intelligence*, 3, 62. <https://doi.org/10.3389/frai.2020.00062>
- [15] Benjamin, R. (2019). *Race after technology: Abolitionist tools for the new Jim Code*. Polity Press.
- [16] Sun, T., Gaut, A., Tang, S., Huang, Y., ElSherief, M., Zhao, J., Mirza, D., Belding, E., Chang, K.-W., & Wang, W. Y. (2019). Mitigating gender bias in natural language processing: Literature review. In *Proceedings of the 57th ACL*, 1630–1640. <https://doi.org/10.18653/v1/P19-1159>
- [17] Lauscher, A., Glavaš, G., Ponzetto, S. P., & Vulić, I. (2020). A little bit is worse than none: Ranking with incomplete information and the design of (better) no-preference options. *Transactions of the ACL*, 8, 519–534. https://doi.org/10.1162/tacl_a_00326
- [18] Shah, D. J., Schwartz, H. A., & Hovy, D. (2020). Predictive biases in natural language processing models: A conceptual framework and overview. In *Proceedings of the 58th ACL*, 5248–5264. <https://doi.org/10.18653/v1/2020.acl-main.468>
- [19] Braun, V., & Clarke, V. (2006). Using thematic analysis in psychology. *Qualitative Research in Psychology*, 3(2), 77–101. <https://doi.org/10.1191/1478088706qp063oa>
- [20] Thomas, J., & Harden, A. (2008). Methods for the thematic synthesis of qualitative research in systematic reviews. *BMC Medical Research Methodology*, 8(1), 45. <https://doi.org/10.1186/1471-2288-8-45>
- [21] Crossley, S. A., Kyle, K., & McNamara, D. S. (2019). Sentiment analysis and social cognition engine (SEANCE): An automatic tool for sentiment, social cognition, and social-order analysis. *Behavior Research Methods*, 49(3), 803–821. <https://doi.org/10.3758/s13428-016-0743-z>
- [22] Koenecke, A., Nam, A., Lake, E., Nudell, J., Quartey, M., Mengesha, Z., Toups, C., Rickford, J. R., Jurafsky, D., & Goel, S. (2020). Racial disparities in automated speech recognition. *Proceedings of the National Academy of Sciences*, 117(14), 7684–7689. <https://doi.org/10.1073/pnas.1915768117>
- [23] Selwyn, N. (2019). *Should robots replace teachers? AI and the future of education*. Polity Press. <https://doi.org/10.2307/j.ctvvc3rbd>