Exploring Dataset Size and Diversity for OCR Post-Correction with hmByT5 Models

Anonymous ACL submission

Abstract

This study explores the application of the 2 hmByT5 model for Optical Character 3 Recognition (OCR) post-correction, 4 focusing on historical German job 5 advertisements. Two versions of the 6 model-standard and fine-tuned on the ICDAR-2019 dataset-were evaluated 8 across subsets of the JobAds dataset. The ۵ effects of dataset size and OCR model 10 diversity on post-correction performance 11 were analyzed. Results show that larger 12 training datasets improve performance, but 13 with diminishing returns, suggesting an 14 optimal balance between annotation effort 15 and model effectiveness. Training on 16 outputs from multiple OCR systems 17 enhances generalization with limited data 18 but may introduce conflicting patterns in 19 larger datasets. Fine-tuning on unrelated 20 datasets, such as ICDAR, reduced 21 performance, underscoring the importance 22 of domain alignment in pre-training. 23

24 1 Introduction

25 Accessing data from historical newspapers in the ²⁶ machine-readable form offers unique opportunities 27 for scholars of many disciplines. Indeed, several 28 present or past projects focused on digitized 29 newspapers, e.g. (Doucet et al., 2020; Ehrmann et 30 al., 2020; Manrique-Gomez et al., 2024) with the 31 aim to develop tools or create large collections. challenges in Optical 32 However, Character 33 Recognition (OCR) still remain far from being 34 solved (Jarlbrink & Snickars, 2017; Late & 35 Kumpulainen, 2021; Torget, 2023; Wevers, 2023) 36 and pose problems to keyword search or further 37 automated processing, as OCRed text often 38 contains 'character recognition mistakes, 39 formatting issues, and hyphenation problems' 40 (Guan & Greene, 2024).

Through the process of post-correction, the OCRed text quality can be significantly improved. While a lot of work has been done in this regard and the progress and interest of the community can be documented e.g. by the ICDAR competitions (Rigaud et al., 2019; Chiron et al., 2017), many problems still remain unsolved and in specific languages and domains, custom post-correction models have to be trained or fine-tuned.

While the influence of the amount of training 50 51 data on the model performance was evaluated on 52 synthetical data (Guan & Greene, 2024), finding 53 the least amount of annotation necessary to get a 54 good model is a pressing question. Also, 55 generalization of the post-correction models 56 remains a question. Exploring performance of post-57 correction models trained on data coming from 58 various OCR models would complement studies 59 like (Dannélls & Persson, 2020) who explore the 60 role of an OCR system for post-correction on 61 Swedish data using as SVM model and find that in 62 most cases, the performance of the post-correction 63 decreases when the model is trained on mixed data 64 from different OCR systems.

In this paper, we present a series of evaluations with a fine-tuned Byte-to-Byte Text-to-Text Transfer Transformer (ByT5) model for historical German on the corpus of job advertisements from digitized newspapers. We focus on how the number of annotations affects the performance of the model. We also compare the effect of different OCR models used and whether a model trained on outputs of more OCR models reaches better generalization ability.

75 2 Related Work

⁷⁶ In a survey, Nguyen et al. (2022)distinguish
⁷⁷ between manual approaches that mostly benefit
⁷⁸ from crowdsourcing, and (semi-)automatic
⁷⁹ approaches. The latter are further divided into

⁸⁰ approaches working with isolated words, e.g. using $_{132}$ 3 ⁸¹ a dictionary or merging outputs from several OCR 82 systems, and context-dependent, that benefit from 133 In our experiments, we use two datasets: Part of the ⁸³ information gained from the context. These include ¹³⁴ ICDAR 2019-POCR (Rigaud et al., 2019) and the ⁸⁴ language models, feature-based machine learning ¹³⁵ JobAds dataset created in the course of our project. 85 (ML) models, and sequence-to-sequence models 136 ⁸⁶ which conceive post-correction as a machine- ¹³⁷ (Ströbel & Clematide, 2019) DE2, DE3 and DE7 ⁸⁷ translation task.

89 face challenges when containing proper nouns, 140 written in the Fraktur font in the form of images ⁹⁰ words of specific domains, historical orthographic ¹⁴¹ together with their textual ground-truth, as well as variations (Nguyen et al., 2022), but also ¹⁴² in the early modern latin font. We used 7896 files ⁹² abbreviation, which is the case of our data, we ¹⁴³ with an average length of 1569 (±666) bytes from ⁹³ further focus only on the contemporary approaches ¹⁴⁴ those written in the Fraktur font. This dataset is 94 considering context.

⁹⁶ synthetic data for post-OCR corrections, including ¹⁴⁷ data can reduce the ground truth needed to create a 97 based on glyph similarity. Using dataset of 8 ⁹⁸ languages and several models, including mT5 (Xue ¹⁴⁹ data shows a SacreBLEU (Post, 2018) score of 99 et al., 2020), mBART (Tang et al., 2020) or BvT5 100 (Xue et al., 2022), they explore the impact of data ¹⁵¹ 101 volume and various methods of data generation on ¹⁵² (Österreichische Nationalbibliothek, 2021), a 102 the model's performance. Finding that ByT5 153 collection of digitized historical newspapers ¹⁰³ performs the best for the post-OCR correction, they ¹⁵⁴ predominantly in German. From 29 different 104 also find that for all languages, the results are the 155 newspaper titles, we picked issues from 1850-105 best when using the whole dataset instead of its 156 1950. From this corpus, a random sample of 1 page fraction for training, and that the data augmentation ¹⁵⁷ per year per newspaper was made. On these pages, 106 x4 reaches generally the best improvement.

108 109 results in (Maheshwari et al., 2022) on Sanskrit or 160 annotated ads were OCRed with the frak2021 110 in (Löfgren & Dannélls, 2024) for Swedish 161 model (Mannheim University Library, 2021) and ¹¹¹ newspapers, encouraging further use of this model. ¹⁶² manually corrected using the Transkribus platform ¹¹² On the other hand, (Debaene et al., 2025) report ¹⁶³ (Kahle et al., 2017). This resulted in 9680 machine-113 that in the post-correction of early modern Dutch ¹⁶⁴ readable, proof-read job advertisements consisting 114 theatre plays, the ByT5 model unnecessarily 165 of on average 196 (±159) bytes. Additional OCR modifies sequences of the OCRed text. 115

116 117 OCR correction and demonstrated its usefulness ¹⁶⁸ 2024) 118 for this task. Later, Thomas et al. (2024) compare 169 University Library, 2023) models. As frak2021 is 119 Llama and BART models on the task of post- 170 the main model used in our project, the evaluation 120 correction of historical newspapers from the 19th 171 dataset contains only OCR text from this model. 121 century. This 122 abbreviations and spelling variations over time. ¹⁷³ 68,46. significantly 174 report that Llama 123 Authors 124 outperforms BART, however, as it was trained 175 encoding an entire job ad is computationally 125 predominantly on English data, its adaptation to 176 infeasible. Therefore, both datasets are split into 126 other languages might prove a challenge. They also 177 segments containing a maximum of 150 bytes. 127 compare the performance with the amount of 128 training data and show that while BART improves 129 significantly with the increased amount of training 130 data, Llama performs very well already from the 131 outset.

Dataset

From the ICDAR 2019-POCR dataset, the DE1 138 (Springmann et al., 2018) were used. They include As the approaches using isolated words often 139 frontpages of newspapers and literary works 145 used to create a fine-tuned hmByt5 model to Guan and Greene (2024) explored generating ¹⁴⁶ investigate whether fine-tuning on unrelated OCR 148 suitable OCR post-correction model. The ICDAR 150 11,87. The fine-tuned model raises this to 79.39.

JobAds dataset was created from ANNO corpus 158 all job advertisements were manually annotated A ByT5 model was also used with very good 159 using doccano (Nakayama et al., 2018). The 166 text for each job ad was created with the Soper et al. (2021) fine-tuned BART for post- 167 german_print (Mannheim University Library, austrian newspaper and (Mannheim corpus is characterized by ¹⁷² The evaluation dataset has a SacreBLEU score of

Because ByT5 encodes text on the byte level,

Methods 178 4

hmByT5 179 **4.1**

180 ByT5 is a variant of the Text-to-Text Transfer 181 Transformer (T5) (Raffel et al., 2020) model

OCR	Model	% of training set				
Variant		10%	25%	50%	75%	100%
Single	Base	65.18	73.39	73.72	74.20	74.60
	ICDAR	65.14	70.36	71.75	72.29	72.96
Multi	Base	72.09	73.31	74.05	74.22	74.42
	ICDAR	68.82	70.91	72.15	72.54	7317

Table 1: SacreBLEU scores of the single-OCR and multi-OCR training variants on the base hmByT5 model and the model fine-tuned on the ICDAR dataset.

182 designed to process text at the byte level rather than 223 optimizer for training, as this was used in the 183 relying on subword tokenization. This architecture 224 original ByT5 model. Training is conducted for 184 enables ByT5 to handle diverse languages, 225 five epochs and evaluated periodically with 185 character sets, and noisy text more effectively. By 226 SacreBLEU. 186 operating directly on raw byte sequences, the 187 model avoids tokenization biases and can better 227 5 188 capture fine-grained character-level patterns. 189 While this approach increases computational 228 In this section, we are presenting SacreBLEU 190 complexity, it enhances robustness in tasks 229 scores for the two models trained on different 191 involving text with unconventional structures, ¹⁹² misspellings, or rare word forms. ByT5 has ¹⁹³ demonstrated strong performance across a variety 194 of natural language processing applications, 195 particularly in scenarios where traditional tokenization methods struggle.

hmByT5¹ are a set of multi-lingual Byt5 models 197 198 fine-tuned on historical data in English, German, 199 French, Finnish, Swedish and Dutch. Here, we use 200 byt5-small-historic-multilingual-span20-flax as it 238 The results of this study reveal several notable

²⁰¹ shows the best performance on German.

202 4.2 Training

204 different versions of the hmByt5 model for OCR 243 and domain adaptation. Additionally, the relatively 205 post-correction: The standard hmByt5 model and a 244 lower performance of the fine-tuned model raises 206 fine-tuned version trained on a subset of the 245 important questions about the suitability of pre-207 ICDAR-2019 dataset. The training data for our 246 training datasets and their impact on downstream 208 experiments is derived from the Anno-dataset, 247 tasks. 209 which consists of historical OCR outputs and their 248 210 corresponding corrected versions by creating a 249 scores with increasing dataset size are consistent 211 randomized 80-20 split. To assess the impact of 250 with established findings in machine learning, 212 dataset size on model performance, we create five 251 where larger training datasets enable models to 213 different subsets containing 10%, 25%, 50%, 75%, 252 capture more diverse patterns and relationships. 214 and 100% of the available training data. 253 The most significant performance gain occurs 215 Additionally, we introduce two variations for each 254 between using 10% and 25% of the training data. subset: A single-OCR variant, which consists of 255 Beyond this point, further improvements become 217 OCR text produced exclusively by the frak2021 256 more moderate across all model configurations. 218 model and a multi-OCR variant, which contains 257 For example, the base hmByT5 model trained on 219 text generated by the three OCR models mentioned 258 the single-OCR variant achieves a 7.2% increase in in the previous section. 220

221 222 dataset configurations. We use the Adafactor 261 total improvement of 9.0%. Considering the

Results

230 subsets of our training data. Table 1: SacreBLEU 231 scores of the single-OCR and multi-OCR training 232 variants on the base hmByT5 model and the model 233 fine-tuned on the ICDAR dataset. shows how well 234 each model performed for five different dataset 235 sizes when trained on either single- or multi-OCR 236 versions of the JobAds data.

237 6 Discussion

²³⁹ trends and provide valuable insights into the factors 240 influencing OCR post-correction performance. The 241 analysis of dataset size and OCR model diversity 203 In this study, we investigate the performance of two 242 highlights critical aspects of model generalization

The observed improvements in SacreBLEU 259 SacreBLEU when moving from 10% to 25% of the Each hmByt5 variant is trained separately on all 260 training data, while using the entire dataset yields a

¹ https://github.com/stefan-it/hmByT5

²⁶² substantial manual effort required to create four ³¹³ can introduce conflicting patterns in larger datasets. 263 times as much training data, the trade-off of slightly 314 Fine-tuning on the ICDAR dataset consistently ²⁶⁴ lower performance might be acceptable in real- ³¹⁵ reduced performance, likely due to domain world applications. 265

Fine-tuning the hmByT5 model on the ICDAR 317 importance of dataset alignment in pre-training. dataset consistently lowers performance across all 318 267 266 training sizes. Despite the ICDAR dataset's focus 319 design and training strategies. Future work could 269 on OCR error correction, it was created with older 320 explore domain-adaptive fine-tuning, synthetic OCR models and features lower OCR quality. This 321 data augmentation, and more nuanced evaluation suggests that error correction does not generalize 322 metrics to develop robust OCR post-correction 271 272 well between datasets with significantly different 323 systems, improving access to digitized historical 273 qualities or originating from different OCR 324 texts. 274 models. Another key factor may be the ICDAR $_{275}$ dataset's larger size relative to the JobAds dataset, $_{325}$ 8 276 as it consists of longer texts. This likely results in 277 the model overfitting to the specific error patterns ³²⁶ A significant limitation of this study is its focus in the ICDAR dataset. This overfitting could ³²⁷ solely on correcting errors in texts that have already 279 explain why SacreBLEU scores for the fine-tuned ³²⁸ undergone OCR processing. The performance of 280 model exhibit higher absolute changes as 329 OCR models is heavily influenced by the preceding additional training data is introduced, progressively ³³⁰ layout analysis step, which determines the structure 282 mitigating this overtuned correction behavior. ²⁸³ Notably, this effect is further supported by the fact ³³² extraction (Liebl & Burghardt, 2021). The 284 that the multi-OCR training variant has a more ³³³ generalizability of post-correction methods may 285 pronounced impact on the fine-tuned ICDAR 286 model, highlighting its sensitivity to training 287 diversity.

In contrast, the multi-OCR training variant 288 produces mixed results for the base hmByT5 290 model. When using only 10% of the training data, the model's performance shifts from a 4.8% 291 reduction in SacreBLEU to a 5.3% improvement. 292 However, when trained on the entire dataset, the single-OCR variant. This suggests that while multi-295 OCR training can improve generalizability when ²⁹⁷ data is limited, it may lead to reduced performance ³⁴⁶ A more extensive dataset would allow for a deeper ²⁹⁸ on larger datasets, possibly due to conflicting error ³⁴⁷ exploration of these trends and provide more robust patterns among different OCR systems. This is also 299 300 in agreement with older findings comparing post-301 correction performance across different OCR 302 systems (Dannélls & Persson, 2020).

Conclusion 303 7

304 This study highlights key factors influencing OCR 305 post-correction performance, including dataset 306 size, OCR model diversity, and pre-training 307 strategies. Larger training datasets improve model 308 performance, though gains diminish as dataset size 309 increases, emphasizing the need for balancing annotation effort and performance benefits. 310 Training on multi-OCR outputs enhances ³⁵⁹ (2017). ICDAR2017 Competition on Post-OCR Text

312 generalization, especially with smaller datasets, but ³⁶⁰ Correction. 2017 14th IAPR International Conference

316 mismatch and overfitting, underscoring the

These findings stress the need for careful dataset

Limitations

331 and organization of the document before text 334 vary depending on the specific techniques 335 employed during layout analysis, even when the 336 same OCR model is used. This aspect warrants 337 further investigation to better understand its impact 338 on OCR post-correction.

Another limitation lies in the size of the JobAds 339 340 dataset used in this study. While the findings 341 illustrate the relationship between dataset size and 342 improvements in OCR quality through postmulti-OCR variant underperforms compared to the ³⁴³ correction, the observed performance decline of 344 multi-OCR models with increasing amounts of 345 training data requires validation on a larger dataset. 348 conclusions about the scalability of multi-OCR 349 training approaches.

350 Acknowledgments

351

352 References

- 353 C. Rigaud, A. Doucet, M. Coustaty, & J. -P. Moreux.
- 354 (2019). ICDAR 2019 Competition on Post-OCR Text
- 355 Correction. 2019 International Conference on
- 356 Document Analysis and Recognition (ICDAR), 1588-
- 357 1593. https://doi.org/10.1109/ICDAR.2019.00255
- 358 Chiron, G., Doucet, A., Coustaty, M., & Moreux, J.-P.

³⁶¹ on Document Analysis and Recognition (ICDAR),	412 Liebl, B., & Burghardt, M. (2021). An Evaluation of
362 1425–1428. https://doi.org/10.1109/ICDAR.2017.252	413 DNN Architectures for Page Segmentation of 414 Historical Newspapers, 2020 25th International
³⁶³ Dannélls, D., & Persson, S. (2020). Supervised OCR	415 Conference on Pattern Recognition (ICPR), 5153–
³⁶⁴ Post-Correction of Historical Swedish Texts: What	416 5160.
366 in the Nordic and Baltic Countries Publications, 3(1),	417 https://doi.org/10.1109/ICPR48806.2021.9412571
367 24–37. https://doi.org/10.5617/dhnbpub.11176	418 Löfgren, V., & Dannélls, D. (2024). Post-OCR
Debeene E. Meledry, A. Lefever E. & Hoste V	419 Correction of Digitized Swedish Newspapers with
(2025) Evaluating Transformers for OCR Post-	420 ByT5. In Y. Bizzoni, S. Degaetano-Ortlieb, A.
370 Correction in Early Modern Dutch Theatre. In O.	421 Kazantseva, & S. Szpakowicz (Hrsg.), Proceedings of
³⁷¹ Rambow, L. Wanner, M. Apidianaki, H. Al-Khalifa,	422 the Sin Joint SIGHOM Workshop on Computational 423 Linguistics for Cultural Heritage Social Sciences
372 B. D. Eugenio, & S. Schockaert (Hrsg.), Proceedings	424 Humanities and Literature (LaTeCH-CLfL 2024) (S.
373 of the 31st International Conference on	425 237–242). Association for Computational Linguistics.
374 Computational Linguistics (S. 1036/–103/4).	426 https://aclanthology.org/2024.latechclfl-1.23/
³⁷⁵ Association for Computational Englishes.	427 Maheshwari A Singh N Krishna A &
sio haps://www.antiorogy.org/2020.com/g main/090/	⁴²⁸ Ramakrishnan, G. (2022). A Benchmark and Dataset
377 Doucet, A., Gasteiner, M., Granroth-Wilding, M.,	429 for Post-OCR text correction in Sanskrit. Findings of
378 Kaiser, M., Kaukonen, M., Labahn, R., Moreux, JP.,	430 the Association for Computational Linguistics:
³⁷⁹ Muemberger, G., Flanzener, E., Therenty, ME., ³⁸⁰ Toivonen H & Tolonen M (2020 Juni) NewsEve:	431 <i>EMNLP</i> 2022, 6258–6265.
³⁶⁰ Forvonen, H., & Foronen, H. (2020, Juli). <i>NewsDyc.</i> ³⁸¹ A digital investigator for historical newspapers.	432 https://doi.org/10.18653/v1/2022.findings-emnlp.466
382 Digital Humanities 2020 (DH 2020), Ottawa, Canada,	433 Mannheim University Library. (2021). Frak2021
383 20-25 July 2020.	434 (Version frak2021-0.905) [Software]. https://ub-
384 https://doi.org/10.5281/zenodo.3895269	435 backup.bib.uni-
³⁸⁵ Ehrmann, M., Romanello, M., Clematide, S., Ströbel,	436 manneim.de/~stweil/tesstrain/Irak2021/tessdata_best
386 P. B., & Barman, R. (2020). Language Resources for	437 / Hak2021-0.703.traineddata
³⁸⁷ Historical Newspapers: The Impresso Collection.	438 Mannheim University Library. (2023).
³⁸⁸ Proceedings of the 12th Language Resources and Evaluation Conference, 058, 068	439 Austrian_newspapers [Software]. https://ub-
³⁸⁹ Evaluation Conference, 936–968. ³⁹⁰ https://doi.org/10.5167/uzh-191270	440 backup.blb.um-
	442 papers/
³⁹¹ Guan, S., & Greene, D. (2024). Advancing Post-OCR	Manuhaina Ulainanita Lihanna (2024). Camana aniat
³⁹² Correction: A Comparative Study of Synthetic Data.	443 Mannaelm University Library. (2024). German_print
³³⁴ Linguistics ACL 2024, 6036–6047.	445 mannheim.de/~stweil/tesstrain/german print/
395 https://doi.org/10.18653/v1/2024.findings-acl.361	
an Iarlbrink I & Snickara P (2017) Cultural baritage	446 Manrique-Gomez, L., Montes, I., Rodriguez Herrera,
as digital noise: Nineteenth century newspapers in the	447 A., & Mainique, R. (2024). Historical Inc. 1941 448 Century Latin American Spanish Newspaper Corpus
³⁹⁸ digital archive. <i>Journal of Documentation</i> , 73(6),	449 with LLM OCR Correction. In M. Hämäläinen, E.
³⁹⁹ 1228–1243. https://doi.org/10.1108/JD-09-2016-0106	450 Öhman, S. Miyagawa, K. Alnajjar, & Y. Bizzoni
400 Kahle P. Colluto S. Hackl. G. & Mühlberger, G.	451 (Hrsg.), Proceedings of the 4th International
401 (2017). Transkribus—A Service Platform for	452 Conference on Natural Language Processing for
402 Transcription, Recognition and Retrieval of Historical	453 Digital Humanities (S. 152–159). Association for 454 Computational Linguistics
403 Documents. 2017 14th IAPR International	455 https://aclanthology.org/2024.nlp4dh-1.13
404 Conference on Document Analysis and Recognition	
405 (ICDAR), 04, 19–24. 405 https://doi.org/10.1109/ICDAR.2017.307	456 Nakayama, H., Kubo, I., Kamura, J., Taniguchi, Y.,
100 mpo.//doi.org/10.110//1CD/mc.201/.30/	458 for Human. https://github.com/doccano/doccano
⁴⁰⁷ Late, E., & Kumpulainen, S. (2021). Interacting with	
408 digitised historical newspapers: Understanding the use	459 Nguyen, I. I. H., Jatowt, A., Coustaty, M., & Doucet, 460 A (2022) Survey of Post OCP Processing
410 Documentation, 78(7). 106–124.	461 Approaches, ACM Computing Surveys, 54(6), 1–37.
411 https://doi.org/10.1108/JD-04-2021-0078	462 https://doi.org/10.1145/3453476
1 8	

- ⁴⁶³ Österreichische Nationalbibliothek. (2021). ANNO
- 464 Historische Zeitungen und Zeitschriften.
- 465 https://anno.onb.ac.at/
- ⁴⁶⁶ Post, M. (2018). A Call for Clarity in Reporting
- 467 BLEU Scores. Proceedings of the Third Conference
- 468 on Machine Translation: Research Papers, 186–191.
- 469 https://doi.org/10.18653/v1/W18-6319
- 470 Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang,
- 471 S., Matena, M., Zhou, Y., Li, W., & Liu, P. J. (2020).
- 472 Exploring the limits of transfer learning with a unified
- 473 text-to-text transformer. J. Mach. Learn. Res., 21(1),
- 474 140:5485-140:5551.
- 475 Soper, E., Fujimoto, S., & Yu, Y.-Y. (2021). BART
- 476 for Post-Correction of OCR Newspaper Text.
- 477 Proceedings of the Seventh Workshop on Noisy User-
- 478 Generated Text (W-NUT 2021), 284–290.
- 479 https://doi.org/10.18653/v1/2021.wnut-1.31
- 480 Springmann, U., Reul, C., Dipper, S., & Baiter, J.
- 481 (2018). Ground Truth for training OCR engines on
- 482 historical documents in German Fraktur and Early
- 483 Modern Latin. arXiv.
- 484 https://doi.org/10.48550/ARXIV.1809.05501
- 485 Ströbel, P., & Clematide, S. (2019). Ground truth for
- ⁴⁸⁶ *Neue Zürcher Zeitung black letter period* [Dataset].
- 487 figshare.
- 488 https://doi.org/10.6084/M9.FIGSHARE.8864510.V1
- 489 Tang, Y., Tran, C., Li, X., Chen, P.-J., Goyal, N.,
- 490 Chaudhary, V., Gu, J., & Fan, A. (2020). Multilingual
- 491 Translation with Extensible Multilingual Pretraining
- 492 and Finetuning. arXiv.
- 493 https://doi.org/10.48550/ARXIV.2008.00401
- 494 Thomas, A., Gaizauskas, R., & Lu, H. (2024).
- 495 Leveraging LLMs for Post-OCR Correction of
- 496 Historical Newspapers. In R. Sprugnoli & M.
- 497 Passarotti (Hrsg.), Proceedings of the Third Workshop
- 498 on Language Technologies for Historical and Ancient
- 499 Languages (LT4HALA) @ LREC-COLING-2024 (S.
- 500 116-121). ELRA and ICCL.
- 501 https://aclanthology.org/2024.lt4hala-1.14/
- 502 Torget, A. (2023). Mapping Texts: Examining the
- 503 Effects of OCR Noise on Historical Newspaper
- 504 Collections. In E. Bunout, M. Ehrmann, & F. Clavert
- 505 (Hrsg.), Digitised Newspapers A New Eldorado for
- 506 Historians? (S. 47-66).
- 507 https://doi.org/10.1515/9783110729214-003
- 508 Wevers, M. (2023). Mining Historical Advertisements
- 509 in Digitised Newspapers. In E. Bunout, M. Ehrmann,
- 510 & F. Clavert (Hrsg.), *Digitised Newspapers A New*
- 511 Eldorado for Historians? (S. 227-252).
- 512 https://doi.org/10.1515/9783110729214-011

- 513 Xue, L., Barua, A., Constant, N., Al-Rfou, R.,
- 514 Narang, S., Kale, M., Roberts, A., & Raffel, C.
- 515 (2022). ByT5: Towards a Token-Free Future with
- 516 Pre-trained Byte-to-Byte Models. Transactions of the
- 517 Association for Computational Linguistics, 10, 291–
- 518 306. https://doi.org/10.1162/tacl_a_00461
- 519 Xue, L., Constant, N., Roberts, A., Kale, M., Al-Rfou,
- 520 R., Siddhant, A., Barua, A., & Raffel, C. (2020). mT5:
- 521 A massively multilingual pre-trained text-to-text
- 522 transformer. arXiv.
- 523 https://doi.org/10.48550/ARXIV.2010.11934
- 524