

Pipeline Stage	Model	Statistic	Result	Used
OCR	EffOCR-Word (Small) (3)	CER	0.015	*
	Tesseract OCR	CER	0.106	*
	EasyOCR	CER	0.170	
	PaddleOCR	CER	0.304	
Deduplication	Locally Sensitive Hashing	ARI	73.7	
	N-gram Overlap	ARI	75.0	
	Neural biencoder (7)	ARI	91.5	*
Georeferencing	Biencoder + crossencoder	ARI	93.7	
	N-gram matching	Accuracy	94.9	*
	GPT-4o-mini	Accuracy	85.3	
NER	Custom NER (6)	F1	90.4	*
	Roberta-Large tuned on CoNLL03 (4)	F1	77.8	
Entity disambiguation	LinkNewsWikipedia (1)	Accuracy	78.3	*
	BLINK (8)	Accuracy	59.9	
	GENRE (5)	Accuracy	63.4	
	ReFinED (2)	Accuracy	65.4	

Table 1: Models considered for each stage of the pipeline and comparisons of their performance on the test set for each stage. Starred models are those used. For CER, smaller values are better, while for the other statistics, bigger numbers are better. CER = Character Error Rate, ARI = Adjusted Rand Index.

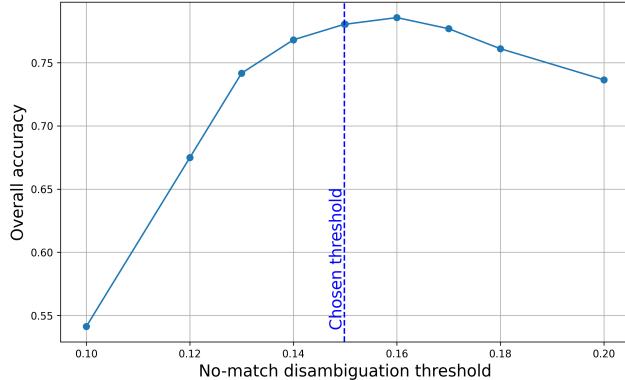


Figure 1: Sensitivity of disambiguation results to choice of no-match threshold.

1 References

- 2 [1] ARORA, A., YANG, X., JHENG, S. Y., AND DELL, M. Linking representations with multimodal contrastive learning. *arXiv preprint arXiv:2304.03464* (2023).
- 3
- 4 [2] AYOOLA, T., TYAGI, S., FISHER, J., CHRISTODOULOUPOULOS, C., AND PIERLEONI, A. ReFinED: An efficient zero-shot-capable
- 5 approach to end-to-end entity linking. In *NAACL* (2022).
- 6 [3] BRYAN, T., CARLSON, J., ARORA, A., AND DELL, M. Efficientocr: An extensible, open-source package for efficiently digitizing world
- 7 knowledge". *Empirical Methods on Natural Language Processing (Systems Demonstrations Track)* (2023).
- 8 [4] CONNEAU, A., KHANDELWAL, K., GOYAL, N., CHAUDHARY, V., WENZEK, G., GUZMÁN, F., GRAVE, E., OTT, M., ZETTLEMOYER,
- 9 L., AND STOYANOV, V. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116* (2019).
- 10 [5] DE CAO, N., IZACARD, G., RIEDEL, S., AND PETRONI, F. Autoregressive entity retrieval. *arXiv preprint arXiv:2010.00904* (2020).
- 11 [6] FRANKLIN, B., SILCOCK, E., ARORA, A., BRYAN, T., AND DELL, M. News dejá vu: Connecting past and present with semantic search,
- 12 2024.
- 13 [7] SILCOCK, E., D'AMICO-WONG, L., YANG, J., AND DELL, M. Noise-robust de-duplication at scale. In *The Eleventh International Conference on Learning Representations* (2022).
- 14
- 15 [8] WU, L., PETRONI, F., JOSIFOSKI, M., RIEDEL, S., AND ZETTLEMOYER, L. Scalable zero-shot entity linking with dense entity retrieval.
- 16 *arXiv preprint arXiv:1911.03814* (2019).