

Yes but.. Can ChatGPT Identify Entities in Historical Documents?

Carlos-Emiliano
González-Gallardo
carlos.gonzalez_gallardo@univ-lr.fr
University of La Rochelle, L3i
La Rochelle, France

Emanuela Boros*
emanuela.boros@epfl.ch
Digital Humanities Laboratory, EPFL
Vaud, Switzerland

Nancy Girdhar
nancy.girdhar@univ-lr.fr
University of La Rochelle, L3i
La Rochelle, France

Ahmed Hamdi
ahmed.hamdi@univ-lr.fr
University of La Rochelle, L3i
La Rochelle, France

Jose G. Moreno
jose.moreno@irit.fr
University of Toulouse, IRIT
Toulouse, France

Antoine Doucet
antoine.doucet@univ-lr.fr
University of La Rochelle, L3i
La Rochelle, France

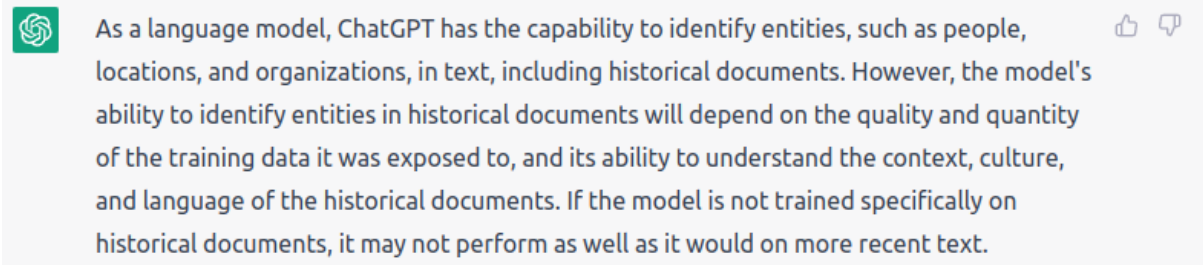


Figure 1: Asking ChatGPT to identify named entities in a historical document.

ABSTRACT

Large language models (LLMs) have been leveraged for several years now, obtaining state-of-the-art performance in recognizing entities from modern documents. For the last few months, the conversational agent ChatGPT has “prompted” a lot of interest in the scientific community and public due to its capacity of generating plausible-sounding answers. In this paper, we explore this ability by probing it in the named entity recognition and classification (NERC) task in primary sources (e.g., historical newspapers and classical commentaries) in a zero-shot manner and by comparing it with state-of-the-art LM-based systems. Our findings indicate several shortcomings in identifying entities in historical text that range from the consistency of entity annotation guidelines, entity complexity, and code-switching, to the specificity of prompting. Moreover, as expected, the inaccessibility of historical archives to the public (and thus on the Internet) also impacts its performance.

CCS CONCEPTS

• **Information systems** → **Language models; Information extraction**; • **Computing methodologies** → **Natural language processing**; • **Applied computing** → **Arts and humanities**.

KEYWORDS

Named entity recognition and classification, Large language models, Generative pretrained transformer, Historical documents

1 INTRODUCTION

Since OpenAI¹ released ChatGPT at the thirty-sixth conference on neural information processing systems (NeurIPS) in November 2022², its ability to provide human-like and plausible-sounding answers caused the model to become extremely popular beyond the research community, gaining more than 1 million users in less than one week. ChatGPT is a conversational agent based on GPT-3.5 (generative pretrained transformer), a large language model (LLM) with more than 175 billion parameters [23]. Given its widespread popularity and accessibility, the question of how this highly mediated model performs on different natural language processing (NLP) tasks arose already in several fields [3, 24].

LLMs have been leveraged for several years now, obtaining state-of-the-art performance in the majority of NLP tasks, by generally being fine-tuned on downstream tasks such as named entity recognition and classification (NERC) and less in zero-shot settings [21]. Thus, for NERC, but also as a general focus, efforts are dedicated to how to effectively transfer knowledge for domain adaptation by developing cross-domain robust systems and exploring zero-shot, or few-shot learning to address domain and annotation consistency and mismatch in cross-domain settings [14, 15]. Simultaneously, in historical documents (e.g., historical newspapers and broadcasts), NERC faces new challenges apart from domain heterogeneity such as input noisiness, language dynamics and lack of resources [12].

Processes such as optical character recognition (OCR) or optical layout recognition (OLR) affect consistently the transcriptions and thus, this propagated noise influences the precision of NERC systems [4, 5, 18, 22, 29, 30]. Even though the latest developments

*This work was done while at University of La Rochelle, in La Rochelle, France.

¹<https://openai.com/>

²<https://nips.cc/>

Table 1: Dataset prompts used to collect the predictions.

NewsEye	hipe-2020	ajmc
What are the locations (LOC), persons (PER), organizations (ORG) and human productions (HumanProd) present in the following historical text? <i>{SENTENCE}</i> Respond, for each word, using IOB or BIO format separated by tab. If a word has no entity, add O.	What are the locations (loc), persons (pers), organizations (org), products (prod) and time periods (time) present in the following historical text? <i>{SENTENCE}</i> Respond, for each word, using IOB or BIO format separated by tab. If a word has no entity, add O.	What are the locations (loc), persons (pers), time periods (date), human works (work), physical objects (object) and specific portion of works (scope) present in the following historical text? <i>{SENTENCE}</i> Respond, for each word, using IOB or BIO format separated by tab. If a word has no entity, add O.

in deep learning by fine-tuning and pre-training historical LMs brought state-of-the-art results in NERC in historical documents [2, 17], time and domain shifts and resource scarcity remain crucial challenges in learning or reusing appropriate knowledge for NERC. Thus, as expected, LLMs and systems in which they are embedded such as ChatGPT were not explicitly trained for information extraction tasks [26] (e.g., named entity recognition and classification, relation extraction), and moreover, as seen in Figure 1, not specifically with a focus in historical documents.

In this short preliminary work, we conduct an exploratory case study to investigate the potential of ChatGPT, which was trained on a massive amount of Internet data (e.g., Common Crawl, WebText2, Wikipedia) [7] and prompt datasets for reinforcement learning from human feedback (RLHF) [23]. Due to this increase in scale in comparison with previous generative language models such as GPT-3 and GPT-2 [7, 16], the behaviour of the model drastically changed, being considerably more able to perform tasks it was not explicitly trained on than previous models [25]. We would expect ChatGPT to be able to detect entities in historical documents to a certain degree, considering the aforementioned challenges, and therefore, we conduct this study by experimenting with zero-shot NERC and comparing its performance against state-of-the-art systems.

2 DATASETS

We selected three historical document datasets covering circa 200 years, they are composed of classical commentaries and historical newspapers, provided by digital libraries during different international research projects, i.e., NewsEye³ and *impresso*⁴.

The NewsEye dataset [19] was collected through the national libraries of France (BnF), Austria (ONB) and Finland (NLF)⁵. It comprises four corpora (French, German, Finnish, and Swedish), being the French one composed of texts from digitized archives of nine newspapers (i.e., *L’Oeuvre*, *La Fronde*, *La Presse*, *Le Matin*, *Marie-Claire*, *Ce soir*, *Marianne*, *Paris Soir* and *Regards*) from 1854 to 1946.

The hipe-2020 dataset [13] is composed of Swiss, Luxembourgish and American newspaper articles in French, German and English comprising the 19th and 20th centuries. It has been collected mainly through the National Library of Switzerland (BN), the National Library of Luxembourg (BnL), the Media Center and State

Table 2: Statistical description of corpora (PERS = person, LOC = location, ORG = organization, HumanProd = human work/production, TIME = date/interval, SCOPE = specific portion of work).

Tokens	Entities	PERS	LOC	ORG	HumanProd	TIME	SCOPE
NewsEye							
30,458	1,298	463	597	217	21	-	-
hipe-2020							
48,854	1,600	502	854	130	61	53	-
ajmc							
5,390	360	139	9	-	80	3	129

Archives of Valais and the Swiss Economic Archives (SWA)⁶ as part of the *impresso* project.

The ajmc dataset [28] is composed of classical commentaries from the Ajax Multi-Commentary project that includes digitized 19th century commentaries published in French, German, and English. These commentaries provide in-depth analysis and explanation of Sophocles’ Ajax Greek tragedy.

These datasets have been annotated with universal (i.e., person, location, organization) and domain-specific (i.e., bibliographic references to primary and secondary literature) entity types and subtypes for the NERC task and split into train, development and test partitions. During this preliminary work, only French test partitions have been taken into consideration. Table 2 displays the information regarding the number and type of entities found in the specified datasets.

3 METHODOLOGY

We followed a straightforward zero-shot approach to retrieve named entities from ChatGPT via the official web interface⁷ between January 11th and February 7th, 2023. An upgrade was released on January 30th to improve factuality and mathematical capabilities of the model⁸, however, we did not perceive any difference with regard to the ability of the model to detect entities.

hipe-2020 and ajmc were provided with fine-grained and nested entities, but in order to simplify the complexity of the prompts for

³<https://www.newseye.eu/>

⁴<https://impresso-project.ch/>

⁵BnF: <https://bnf.fr/>; ONB: <https://onb.ac.at/>; NLF: <https://kansalliskirjasto.fi>

⁶BN: <https://www.nb.admin.ch/>; BnL: <https://bnl.public.lu/>; SWA: <https://wirtschaftsarchiv.ub.unibas.ch>

⁷<https://chat.openai.com>

⁸<https://help.openai.com/en/articles/6825453-chatgpt-release-notes>

this study, we only consider the coarse-grained entities. We defined the three prompts presented in Table 1 depending on the entity type differences between datasets and to respect their corresponding labels casing. Even if the IOB/BIO⁹ format is explicitly demanded for each word, tokenization by ChatGPT was inconsistent with IOB tokenized dataset files, thus, for evaluation purposes, an alignment and verification process was required to ensure evaluation consistency.

4 RESULTS

Table 3 presents the performance of ChatGPT over coarse-grained (high-level entity types) NERC in terms of strict and fuzzy micro-level precision (P), recall (R) and F-measure (F1) evaluated with CLEF-HIPE-2020-scorer¹⁰. For positioning and comparison, we also present the performance of two LM-based state-of-the-art NERC systems.

Table 3: Comparative results using the three datasets (micro).

	NewsEye			hipe-2020			ajmc		
	P	R	F1	P	R	F1	P	R	F1
	strict								
<i>Stacked NERC</i>	75.0	70.6	72.7	-	-	-	-	-	-
<i>Temporal NERC</i>	-	-	-	76.5	76.5	76.5	84.8	83.9	84.4
ChatGPT	70.9	72.3	71.6	32.5	50.0	39.4	21.8	26.1	23.8
	fuzzy								
<i>Stacked NERC</i>	85.4	80.5	82.9	-	-	-	-	-	-
<i>Temporal NERC</i>	-	-	-	86.7	86.7	86.7	90.2	89.2	89.7
ChatGPT	77.8	79.4	78.6	49.0	75.4	59.4	25.5	30.6	27.8

Stacked NERC. This model is based on the pre-trained model BERT proposed by [10] with a stack of 2 Transformer blocks on top, finalized with a conditional random field (CRF) prediction layer. *Stacked NERC* proved increased performance in historical documents, while did not degrade the performance over modern data [5, 6]. The same architecture was utilized as a baseline in the description of the NewsEye dataset [19].

Temporal NERC. This model relies on *Stacked NERC*, and it includes a data-wise improvement by exploiting temporal knowledge graphs for generating additional contextual time information and a model-wise improvement that incorporates this information with mean-pooled context jokers [18]. *Temporal NERC* proved the importance of temporality for historical newspapers and classical commentaries, depending on the time intervals and the digitization error rate.

From Table 3 it is clear that the capacity of ChatGPT of identifying named entities is really dependent on the dataset and the type of entities. Noticeable drops in performance can be observed for ajmc with over 71% decrease for strict and over 69% for fuzzy metrics. For hipe-2020, the performance has decreased less drastically, with over 48% for strict and 31.48% for fuzzy. With respect

⁹[https://en.wikipedia.org/wiki/Inside%E2%80%93outside%E2%80%93beginning_\(tagging\)](https://en.wikipedia.org/wiki/Inside%E2%80%93outside%E2%80%93beginning_(tagging))

¹⁰<https://github.com/hipe-eval/HIPE-scorer>

to NewsEye, the score values are marginally similar, with a drop of close to 1.5% for strict and slightly over 5% for fuzzy. While the results are generally balanced, we also observe a higher recall in the case of hipe-2020 which could indicate that the complexity of the entity annotation causes ChatGPT to be able to detect them but not correctly classify them.

5 LIMITATIONS

We refer to limitations as the weaknesses of which ChatGPT gives evidence in the process of recognizing entities in a historical text, and less regarding limitations of this work which are not in our control, such as the inability of exploring the insights of ChatGPT and ensuring reproducible results if the same method is applied¹¹. We, thus, study the impact of the quality of the datasets and their annotations with regard to the ChatGPT responses.

5.1 Named Entity Definition

Named entity annotation follows well-defined guidelines to describe the nature and boundaries of universal and domain-specific entity types, yet it is necessary to rely on the linguistic intuition and awareness of the annotator [11, 20, 27]. While the definitions of universal entity types are stable (e.g., person, location and organization), domain-specific entity types vary among guidelines.

hipe-2020 [11] and NewsEye [20] guidelines define a “human production” entity as anything that is broadcast in the press, on radio or television such as newspapers, magazines, broadcasts, sales catalogues, etc. (e.g., “Die Zeit”, “Le Figaro”, “Le sept à huit”, “La ferme célébrités”) and exclude media products such as films, TV series, etc., and political, philosophical, religious/sectarian doctrines (e.g., “Der Sozialismus”, “Theravada Buddhismus”, “Le socialism”, “le bouddhisme theravâda”).

Similarly, a “work” entity is described by the ajmc guidelines [27] as an entity denoting a human creation, be it intellectual or artistic, that can be referred to by its title; “A work is a distinct intellectual or artistic creation” (FRBR¹² guidelines) including literary works, religious works, editions of papyrological and epigraphical sources (e.g. “IG2”, “P.Oxy 1.119”), and journals.

All descriptions overlap, but they undoubtedly lead to different annotations specific to each dataset, that could lead to slight confusions in predictions. For example, in *Le sens est donc : « Ai-je parlé trop peu clairement ? » Cf. Antigone, 405 : “Ap’ Evâenda xxi cœpñ XêYw ; Eschyle, Agamemnon, 269 : “H ζορῶς λέγω ;*, an example from ajmc, where the word *Antigone* refers to a human creation, the entity is not detected by ChatGPT. However, since we are exploring the model in a zero-shot manner, and while we agree that consistency of the annotations is a major concern because of the language ambiguity, we can only assume that this misclassification comes more from the lack of specificity in the prompt than the entity definition.

5.2 Entity Complexity

For NewsEye guidelines, a named entity is a real-world object denoting a unique individual with a proper name. Since, historically, a

¹¹All predictions are available at https://anonymous.4open.science/r/NERC_ChatGPT-F687/.

¹²<https://www.oclc.org/research/activities/frbr.html>

person’s name played an influential role in reflecting key attributes of their job or life, the majority of entities have also been annotated including the person’s job title. In the following example: *Beethoven*. — *Par le Quatuor de la fondation Beethoven : MM.A. Géioso 1er violon ; A. Tracol ; 2e violon ; P. Monteux. aito, F. Schnéklud, violon-celle ; César Geloso, pianiste.* César Geloso, the pianist, is considered a “person” entity, however, ChatGPT is unable to detect beyond the mention of the first and last names. If jobs are needed to be detected inside entities, we assume that more information should be added in the prompt.

If the first or last name are in the text, ChatGPT seems to give them priority because it is able to identify the presence of an entity in the text, even though no names are mentioned (e.g., *garçon, infirmière des Hôpitaux*). Thus, in *Ouvrier d’art et artiste peintre cherche petits travaux restauration de tableaux, décoration d’intérieur ou réparations meubles de style, anciens ou modernes, laques, vernis Martin et au tampon. Christophe, 56, rue de la Montagne-Sainte-Geneviève, 5e 43 ans.* the art worker and painter that looks for small jobs is correctly spotted. As well, in *Un homme à grosses moustaches, pensée, sourit soudain imperceptiblement. Mais La jeune femme du premier rang au manteau de vison, jusqu’alors muette.*, the man with a big moustache was also detected, along with the young woman of the first rank.

Letter casing also seems to be challenging. While it is common for the names of organizations and headlines of newspaper articles to be all capitals (e.g., LE SPORT, NOUVELLES BREVES), ChatGPT identified them all as organizations. In the case of demonyms, locations and persons, they were systematically mixed up (e.g., *Carcassonnais, Russe, Mexicains, “Italien, 17 3/4”, “Japon 1899, 73”, “Portugais 3 %, 2 1/4”, “Russe 1906”*). It is unclear, however, if the confusion comes from the fact that they start with capital letters or due to context, as this limitation is commonly found in recent state-of-the-art NERC models.

NewsEye guidelines consider addresses such as *130, rue de la Courselle* and *56, rue de la Montagne-Sainte-Geneviève, 5e* as locations, however, ChatGPT does not seem to capture this fine level of granularity. Since the word “location” defines space in a more semantically generic manner than a specific place or position, an address is referencing the particulars of a place, which, if unspecified in the prompt, ChatGPT is unable to correctly identify.

5.3 Digitization Errors

Quantitatively, ChatGPT identified only 7% of named entities with OCR errors in the *ajmc* dataset, while *Temporal NERC* correctly identified around 40% of noisy entities. *Temporal NERC* recognized named entities with up to 70% of their characters sustained an OCR error (i.e., deletion, insertion, substitution), however, character errors should not exceed 20% on entities to be recognized by ChatGPT. More specifically, ChatGPT showed an inability to recognize named entities with segmentation errors. For instance, in *13659 - 4360.Hxépra. . . . elloug. Ulysse paraît faire allusion à l’amertume des peroles que vient de prononcer À gamemnon ; Agamemnon répond comme si Ulysse avait eu en vue l’amertume de ses propres remontrances, 4866.ÉvOaë’ Tlouet, j’en arriverai là, c’est-à-dire, je mourrai, .Dindorf : Kai’ αὐτὸς ἴξομαι πρὸς τὸ θάπτειν αὐτόν.*, the name of a person *Agamemnon* has been correctly identified, however, *À*

gamemnon has not been recognized. Not surprisingly, when the entity is highly impacted by the OCR process, such as *Beethoven* instead of *Beethoven* in *Xle Quatuor (op. 95) Beethoven*, ChatGPT fails to detect such entities.

Results also showed that ChatGPT can be affected by spelling variations over years, with 78% of named entities from 1958 to 2018 having been correctly recognized. However, this rate drops to half with entities from 1798 to 1948 in the *hipe-2020* dataset.

Finally, due to an abundant amount of noise, such as in *m_’i’ — .:’i -i’ i’ —’i m Nota. Les Avis à insérer dans cette Feuille, qui ne seront pas remis au Bureau le Mardi matin, à neuf-heures au plus tard, seront renvoyés irrévocablement à l’ordinaire prochain . ARTICLES OFFICIELS. k. Le public est informé, qu’au bénéfice d’un gracieux arrêt de la Seigneurie, “et d’une sentence de direction de l’honorable Justice de la Chaux-de-Fonds, le tuteur et les parens des enfans’ de Henri Jean-Maire, se présenteront en cour de Justice de la die Chatùx-de-Fonds, le Mardi _s_?”, ChatGPT “surrenders” with this statement: *The text you provided is not in a coherent language, and it is difficult to understand what it is trying to communicate. The text contains mostly punctuation and special characters with no recognizable entities.**

5.4 Code-switching

Code-switching refers to the phenomenon of alternating between two or more languages within the same sentence, phrase, or single word. The monolingual bias in multilingualism is a type of bias that can occur in language models [8, 31]. Not only, but these LMs are also hardly generalizable to different code-switched languages [1] and pre-trained multilingual models do not necessarily guarantee high-quality representations on code-switching [32].

ajmc presents code-switching between French and Ancient Greek. GPT-3.5 was trained in more than 100 languages, being English over-represented with 93% word count, 1.82% for French and 0.032% for Modern Greek, while Ancient Greek is unrepresented [7]. ChatGPT was trained with further datasets for fine-tuning, human feedback and prompting [23] but similar language distribution is to be expected. Thus, for an example such as *À la marge d’un exemplaire de Sophocle, on lit la traduction suivante de ces deux vers, due à notre ÆRacine : « O mon fils, sois un jour plus beareux que ton père; Da reste avec honneur tu peux lui ressembler»*. ChatGPT “confidently” responds with *I’m sorry, but I’m unable to understand the text you have provided. The text appears to be a mixture of ancient Greek and French, with some references to ancient literature and annotations, which makes it difficult to extract meaningful information. Additionally, the text is not a historical text, but rather a literary text, which also makes it difficult to extract historical entities.*

6 CONCLUSIONS & PERSPECTIVES

On this note, we conclude that ChatGPT encounters several difficulties in recognizing entities in historical documents, that range from the consistency of entity annotation guidelines, entity complexity, multilingualism and code-switching, to the specificity of prompting. Moreover, while an unprecedented amount of historical documents is available in digital format, little is freely available with many historical archives that remain inaccessible to the public (and consequently, on the Internet). For example, primary sources (documents or other items that provide first-hand, eyewitness accounts

of events) such as newspapers and magazine articles (as in the case of the majority of the datasets from this study) are available both online and in the library, but nevertheless, they are generally watermarked or behind a paywall. Consequently, ChatGPT is “unaware” (for now) of such knowledge, which contributes to the degree of confusion of the model with regard to historical documents. Even if these resources become accessible and LLMs-based systems improve their “understanding” capacities in historical documents, their implementation in digital libraries must be taken with caution to prevent biased and out-of-domain responses.

ETHICS STATEMENT

While it can generate plausible-sounding text, the content generated by ChatGPT is not necessarily true. Nevertheless, we consider that it does not concern the task of named entity recognition, as we are not adding any further ethical consideration other than those posed by ChatGPT. We are aware of the intentional stance [9] of terms such as “confidently”, “surrenders”, “understanding” and “unaware” when applied to a conversational agent, however, in this paper, we adopt them to emphasize ChatGPT capacity to interact with a user.

ACKNOWLEDGEMENTS

This work has been supported by the ANNA (2019-1R40226) and TERMITRAD (2020-2019-8510010) projects funded by the Nouvelle-Aquitaine Region, France.

REFERENCES

- [1] Gustavo Aguilar, Sudipta Kar, and Thamar Solorio. 2020. LinCE: A Centralized Benchmark for Linguistic Code-switching Evaluation. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. European Language Resources Association, Marseille, France, 1803–1813. <https://aclanthology.org/2020.lrec-1.223>
- [2] Alessio Palmero Aprosio, Stefano Menini, and Sara Tonelli. 2022. BERToldo, the Historical BERT for Italian. In *Proceedings of the Second Workshop on Language Technologies for Historical and Ancient Languages*. 68–72.
- [3] Som Biswas. 2023. ChatGPT and the Future of Medical Writing. , 223312 pages.
- [4] Emanuela Boros, Carlos-Emiliano González-Gallardo, Edward Giamphy, Ahmed Hamdi, José G Moreno, and Antoine Doucet. 2022. Knowledge-based Contexts for Historical Named Entity Recognition & Linking. *Conference and Labs of the Evaluation Forum (CLEF 2020)* (2022).
- [5] Emanuela Boros, Ahmed Hamdi, Elvys Linhares Pontes, Luis-Adrián Cabrera-Diego, Jose G Moreno, Nicolas Sidere, and Antoine Doucet. 2020. Alleviating digitization errors in named entity recognition for historical documents. In *Proceedings of the 24th conference on computational natural language learning*. 431–441.
- [6] Emanuela Boros, Elvys Linhares Pontes, Luis Adrián Cabrera-Diego, Ahmed Hamdi, José G Moreno, Nicolas Sidere, and Antoine Doucet. 2020. Robust named entity recognition and linking on historical multilingual documents. In *Conference and Labs of the Evaluation Forum (CLEF 2020)*, Vol. 2696. CEUR-WS Working Notes, 1–17.
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [8] Monojit Choudhury and Amit Deshpande. 2021. How Linguistically Fair Are Multilingual Pre-Trained Language Models?. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 12710–12718.
- [9] Daniel Dennett. 2009. Intentional Systems Theory. In *The Oxford Handbook of Philosophy of Mind*. Oxford University Press. <https://doi.org/10.1093/oxfordhb/9780199262618.003.0020>
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [11] Ehrmann, Watter, Romanello, Clematide, and Flückiger. 2020. *Impresso Named Entity Annotation Guidelines*. <https://doi.org/10.5281/zenodo.3604227>
- [12] Maud Ehrmann, Ahmed Hamdi, Elvys Linhares Pontes, Matteo Romanello, and Antoine Doucet. 2021. Named entity recognition and classification on historical documents: A survey. *arXiv preprint arXiv:2109.11406* (2021).
- [13] Maud Ehrmann, Matteo Romanello, Simon Clematide, Phillip Ströbel, and Raphaël Barman. 2020. Language resources for historical newspapers: the impresso collection. (2020).
- [14] Maud Ehrmann, Matteo Romanello, Alex Flückiger, and Simon Clematide. 2020. Extended overview of CLEF HIPE 2020: named entity processing on historical newspapers. In *CEUR Workshop Proceedings*. CEUR-WS.
- [15] Maud Ehrmann, Matteo Romanello, Sven Najem-Meyer, Antoine Doucet, and Simon Clematide. 2022. Extended Overview of HIPE-2022: Named Entity Recognition and Linking in Multilingual Historical Documents. In *Proceedings of the 13th International Conference of the CLEF Association (Lecture Notes in Computer Science)*, Vol. 13390. Springer.
- [16] Luciano Floridi and Massimo Chiriatti. 2020. GPT-3: Its nature, scope, limits, and consequences. *Minds and Machines* 30 (2020), 681–694.
- [17] Lauren Fonteyn and Enrique Manjavacas. 2022. Adapting vs. Pre-training Language Models for Historical Languages. *Journal of Data Mining & Digital Humanities* (2022).
- [18] Carlos-Emiliano González-Gallardo, Emanuela Boros, Edward Giamphy, Ahmed Hamdi, José G Moreno, and Antoine Doucet. 2023. Injecting Temporal-aware Knowledge in Historical Named Entity Recognition. In *Advances in Information Retrieval: 45th European Conference on IR Research, ECIR 2023, Dublin, Ireland, April 2–6, 2023, Proceedings, Part I*. Springer, 65–79.
- [19] Ahmed Hamdi, Elvys Linhares Pontes, Emanuela Boros, Thi Tuyet Hai Nguyen, Günter Hackl, Jose G Moreno, and Antoine Doucet. 2021. A multilingual dataset for named entity recognition, entity linking and stance detection in historical newspapers. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2328–2334.
- [20] Ahmed Hamdi, Elvys Linhares Pontes, and Antoine Doucet. 2021. Annotation Guidelines for Named Entity Recognition, Entity Linking and Stance Detection. <https://doi.org/10.5281/zenodo.4574199>
- [21] Jing Li, Aixin Sun, Jianglei Han, and Chenliang Li. 2020. A survey on deep learning for named entity recognition. *IEEE Transactions on Knowledge and Data Engineering* 34, 1 (2020), 50–70.
- [22] Sven Najem-Meyer and Matteo Romanello. 2022. Page Layout Analysis of Text-heavy Historical Documents: a Comparison of Textual and Visual Approaches. In *Proceedings of the Computational Humanities Research Conference 2022 Antwerp, Belgium, December 12-14, 2022*. 36–54.
- [23] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155* (2022).
- [24] John V Pavlik. 2023. Collaborating With ChatGPT: Considering the Implications of Generative Artificial Intelligence for Journalism and Media Education. *Journalism & Mass Communication Educator* (2023), 10776958221149577.
- [25] Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. True few-shot learning with language models. *Advances in neural information processing systems* 34 (2021), 11054–11070.
- [26] Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is ChatGPT a General-Purpose Natural Language Processing Task Solver? *arXiv preprint arXiv:2302.06476* (2023).
- [27] Matteo Romanello and Sven Najem-Meyer. 2022. *Guidelines for the Annotation of Named Entities in the Domain of Classics*. <https://doi.org/10.5281/zenodo.6368101>
- [28] Matteo Romanello, Sven Najem-Meyer, and Bruce Robertson. 2021. Optical Character Recognition of 19th Century Classical Commentaries: The Current State of Affairs. In *The 6th International Workshop on Historical Document Imaging and Processing (Lausanne, Switzerland) (HIP '21)*. Association for Computing Machinery, New York, NY, USA, 1–6. <https://doi.org/10.1145/3476887.3476911>
- [29] Stefan Schweter and Johannes Baiter. 2019. Towards Robust Named Entity Recognition for Historic German. In *Proceedings of the 4th Workshop on Representation Learning for NLP (Repl4NLP-2019)*. Association for Computational Linguistics, Florence, Italy, 96–103. <https://doi.org/10.18653/v1/W19-4312>
- [30] Stefan Schweter, Luisa März, Katharina Schmid, and Erion Čano. 2022. hmbert: Historical multilingual language models for named entity recognition. *Conference and Labs of the Evaluation Forum (CLEF 2020)* (2022).
- [31] Zeerak Talat, Aurélie Névéol, Stella Biderman, Miruna Clinciu, Manan Dey, Shayne Longpre, Sasha Luccioni, Maraim Masoud, Margaret Mitchell, Dragomir Radev, et al. 2022. You reap what you sow: On the challenges of bias evaluation under multilingual settings. In *Proceedings of BigScience Episode# 5-Workshop on Challenges & Perspectives in Creating Large Language Models*. 26–41.
- [32] Genta Indra Winata, Samuel Cahyawijaya, Zihan Liu, Zhaojiang Lin, Andrea Madotto, and Pascale Fung. 2021. Are Multilingual Models Effective in Code-Switching?. In *Proceedings of the Fifth Workshop on Computational Approaches to Linguistic Code-Switching*. Association for Computational Linguistics, Online, 142–153. <https://doi.org/10.18653/v1/2021.calcs-1.20>