

A Dual-View Analysis of Multiple Languages in Colonial Newspapers

WARNING: This paper shows dataset samples which are racist in nature

Anonymous ACL submission

Abstract

Historical newspapers from the colonial period offer valuable evidence of how racializing language evolved over time. However, there are challenges in studying this type of historical data: 1) Data scarcity: acquiring large, annotated historical datasets is difficult, hindering the possibility of analyzing racialization comprehensively; 2) Digitized materials frequently contain Optical Character Recognition (OCR) errors and other types of noise that complicate text extraction and computational analysis; 3) Colonial newspapers are often multilingual and written in archaic prose, hindering the effectiveness of NLP tools developed for modern, single language texts. This paper addresses these challenges by conducting a dual-view, jointly studying multilingual event extraction and temporal semantic shift tasks. Specifically, we introduce a contextual question answering (CQA) and a visual question answering (VQA) derived from eighteenth- and nineteenth-century colonial newspapers, which aims to capture the idiosyncratic characteristics of this material. Content-wise, we focus on how enslaved people were described by enslavers as well as how they articulated their own condition through QA pairs of newspapers written in Dutch, English-French, and Spanish. Our results show that LLMs are still limited for low-resource VQA tasks. For temporal semantic change, we train temporal word embedding with a compass. The study concludes that racialization is a fluid process of linguistic recalibration where the decline of slavery merely shifted the language of control onto new categories of labor and identity, leaving a "semantic debt" that still shapes language bias today.

1 Introduction

The historical study of racialization examines how racial identities are formed and maintained over time. This topic is central to research in history (Silva Perez, 2025), sociology, and sociolinguistics

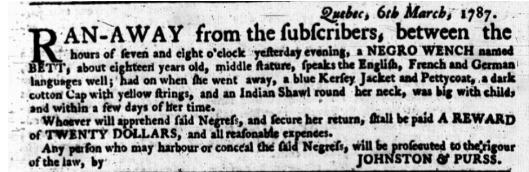


Figure 1: A screenshot from the historical colonial newspapers. The resulting texts often contain unusual spellings, layout noise, and many OCR errors.

(Alim, 2016; Costa and Sokolovska, 2025; Grieve et al., 2025). Within this field, historical colonial newspapers are especially valuable. They served as public forums where ideas of racial hierarchy and resistance were openly expressed (da Silva Perez and Borenstein, 2025; Borenstein et al., 2023a; Newman, 2022; Mäkinen, 2022). Historians traditionally work closely with the texts they study. Nevertheless, NLP tools can automate parts of the analysis and speed up the research process. This allows historians to extract evidence from large corpora more efficiently and focus more on interpretation (Borenstein et al., 2023a).

However, leveraging these archives on NLP tasks remains a profound challenge for researchers (Hill and Hengchen, 2019; Rigaud et al., 2019). There are three main challenges. First, building high-quality datasets still depends heavily on manual work, which is slow and costly. This makes it hard to draw general conclusions about racialized language (Borenstein et al., 2023a). Second, and more importantly, the digitization of these fragile, centuries-old documents creates serious data quality problems. The resulting texts often contain unusual spellings, layout noise, and many OCR errors (Fig. 1). Finally, colonial newspapers are often multilingual and written in archaic languages, hindering the effectiveness of standard NLP tools, which are usually trained on clean, modern text (Hill and Hengchen, 2019; Todorov and Colavizza, 2022).

075 These challenges are especially acute for less- 127
076 studied tasks (Borenstein et al., 2023a,b). In this 128
077 work, we argue that event extraction (EE) and se- 129
078 mantic shift offer two complementary perspectives 130
079 on racialization in colonial newspapers. **On the 131**
080 **one hand**, EE captures the explicit, event-level 132
081 manifestations of racialization, such as acts of 133
082 enslavement, labor control, punishment, or resis- 134
083 tance, as they are described in colonial texts (Sprug- 135
084 noli and Tonelli, 2019; Lai et al., 2021; Dagdelen 136
085 et al., 2024). Prior work shows that extracting 137
086 events from historical advertisements can reveal 138
087 how racialized groups were placed within colonial 139
088 power structures. However, performance remains 140
089 limited because of OCR errors and the use of arch- 141
090 aic language (Borenstein et al., 2023b). **On the 142**
091 **other hand**, the semantic shift task captures the im- 143
092 plicit, long-term evolution of racialized meaning, 144
093 revealing how terms associated with race, labor, 145
094 and identity gradually changed as colonial systems 146
095 transformed (Borenstein et al., 2023b). In sum- 147
096 mary, EE focuses on what happens in specific his- 148
097 torical moments while semantic shift reveals how 149
098 language itself adapts to sustain racial hierarchies 150
099 over time. 151

100 By studying event extraction and semantic shift 152
101 together, we link specific racialized events to 153
102 broader changes in language. This dual view helps 154
103 connect historical actions with evolving meanings 155
104 that still influence modern racial bias. We address 156
105 the shortcomings of previous work and make the 157
106 following contributions: (1) To address the dataset 158
107 scarcity and OCR noisy challenges, we use LLMs 159
108 to study multilingual event extraction in Dutch, 160
109 English-French, and Spanish by constructing con- 161
110 textual question answering (CQA) and visual ques- 162
111 tion answering (VQA), centered on historical racial- 163
112 ization processes in colonial newspapers from 18- 164
113 19 centuries. (2) To study how racializing language 165
114 evolved, we study the semantic shift in historical 166
115 newspapers from the Dutch in the colonial period 167
116 between 1816 and 1882, since the abolition of slav- 168
117 ery in the Dutch colonies in 1863 serves as a use- 169
118 ful case study for examining temporal shifts. We 170
119 train the temporal word embeddings with a com- 171
120 pass (TWEC) framework in each year from 1816 172
121 to 1882 to study the semantic change at the word 173
122 level. 174

123 We find that current LLMs perform well on 175
124 CQA tasks when provided with clean, modern text. 176
125 However, when the context consists of images of 177
126 original colonial newspapers, even recent state-of-

the-art multimodal models perform poorly, high- 127
lighting persistent challenges in processing archaic 128
language and historical print. Beyond model per- 129
formance, our semantic shift analysis shows that 130
racialization is a changing process in language. As 131
formal slavery declined, words of control shifted 132
to new kinds of labor and identity. This created a 133
lasting “semantic debt” that still shapes language 134
bias today¹. 135

2 Related Work 136

2.1 LLM for historical research 137

138 The use of LLMs in historical research has changed 139
139 digital humanities. Research now goes beyond 140
140 keyword search and topic modeling and moves 141
141 toward methods that focus on meaning, genera- 142
142 tion, and active reasoning (Grossmann et al., 2023; 143
143 Piper and Wu, 2025; Cherukuri et al., 2025; Zhang 144
144 et al., 2024b; Simons et al., 2025; Hauser et al., 145
145 2024; Karch et al., 2025; Zhang and Colavizza, 146
146 2025; Levchenko, 2025; Humphries et al., 2024). 147
147 (Levchenko, 2025) demonstrates that multimodal 148
148 LLMs achieved a character error rate (CER) as 149
149 low as 1.8% on 18th and 19th-century English 150
150 manuscripts, surpassing state-of-the-art tools. Be- 151
151 yond direct transcription, LLMs are increasingly 152
152 used to post-correct noisy OCR outputs. (Tudor 153
153 et al., 2025) explored this in the context of "prompt- 154
154 ing the past," using LLMs to clean historical text 155
155 transcripts. (Zhang et al., 2024a) explore the po- 156
156 tential of OpenAI’s GPT models on the post-OCR 157
157 correction task containing English texts published 158
158 between 1559 and 1928 on versification and pro- 159
159 nunciation. Unlike prior studies focusing on us- 160
160 ing LLM to OCR issues, we study the historical 161
161 research by studying the multilingual event extrac- 162
162 tion and tracking the semantic evolution of racial- 163
163 ization. 164

2.2 Event Extraction 164

165 Event Extraction (EE) is a task of organizing natu- 166
166 ral texts into structured events (Hogenboom et al., 167
167 2011; Xiang and Wang, 2019; Chen et al., 2024; 168
168 Ma et al., 2024). Usually, EE is decomposed into 169
169 smaller, less complex tasks (Lin et al., 2020). such 170
170 as detecting the existence of an event (Weng and 171
171 Lee, 2011; Nguyen and Grishman, 2018; Sims 172
172 et al., 2019), identifying its participants (Du et al., 173
173 2021; Li et al., 2020), and extracting the attributes

¹https://anonymous.4open.science/status/historical_nlp-429F

174	associated with the event (Li et al., 2020). Some	and is in Spanish: 50 newspaper issues, ap-	222
175	research work has shown the benefit of framing the	proximately 400 pages, spanning 1878-1882.	223
176	EE as QA tasks (Li et al., 2020). The work most		
177	closely related to ours is (Borenstein et al., 2023a),	We focused on newspapers that contained texts	224
178	which studies event extraction by building multi-	related to slavery, enslaved labour, or racialized	225
179	lingual QA pairs from newspaper advertisements	populations (fugitive ads, sale notices, offers of en-	226
180	in the early modern colonial period. In contrast,	slaved workers for hire, etc), were published in dif-	227
181	their work focuses only on advertisements. We	ferent languages beyond English, in colonies gov-	228
182	use LLMs to generate data from entire newspaper	erned by different empires, and were made avail-	229
183	pages and construct both contextual QA and visual	able in open access by the respective repositories.	230
184	QA tasks.		
185	2.3 Temporal Semantic Change	4 Event Extraction from Colonial	231
186	Historical corpora have been widely used to study	Newspapers	232
187	social phenomena such as language change (Ku-	As discussed in the introduction, we study event	233
188	tuzov et al., 2018; Borenstein et al., 2023b). For	extraction from colonial newspapers by framing it	234
189	example, (Garg et al., 2018) use word embeddings	as a question answering task and using LLMs to	235
190	to measure changes in stereotypes and attitudes to-	extract event arguments. We introduce both CQA	236
191	ward women and ethnic minorities in the United	and VQA benchmarks based on newspapers from	237
192	States during the twentieth and twenty-first cen-	the eighteenth and nineteenth centuries in Dutch,	238
193	turies. (Levis Sullam et al., 2022) analyze a large	French, and Spanish.	239
194	and diverse text collection, including books, news-		
195	papers, songs, and essays, to track how Jews are	4.1 Contextual Question Answering Tasks	240
196	linked to different semantic topics and how these	We formalize the contextual question answering	241
197	links change over time.	task as a tuple (C, q, a) . Here, C denotes the con-	242
198	(Borenstein et al., 2023b) also use word embed-	text, representing a specific historical text segment	243
199	dings to study how bias in historical newspapers	such as a fugitive advertisement or a sales notice,	244
200	persists and changes. However, training separate	which serves as the unstructured source of linguis-	245
201	word embeddings for each year makes it hard to	tic evidence. q represents a query designed to probe	246
202	compare meanings across time. In this paper, we	specific dimensions of racialization, ranging from	247
203	use the TWEC framework to study word-level se-	physical descriptors to markers of resistance, effec-	248
204	mantic change in a shared vector space.	tively acting as a sociological prompt. The objec-	249
205		tive of the model is to map the input pair (C, q) to	250
206	3 Dataset Collection	an answer a , which extracts the relevant semantic	251
207	We constructed our multilingual dataset from three	span or generates a synthesized response grounded	252
208	digitized colonial newspapers available in open ac-	in the text. The dataset preprocessing setting is pre-	253
209		sented in Appendix A.7.2. We evaluate the current	254
210	• <i>De Curaçaosche Courant</i> , held by Delpher.nl,	strong open LLM (e.g., Qwen-72B (Team et al.,	255
211	was published in Willemstad, Curaçao, and is	2024)) and closed-formed LLM (e.g., GPT5 (Hurst	256
212	mostly in Dutch with small portions in other	et al., 2024)). The details are presented in Ap-	257
213	languages: 3249 newspaper issues, approxi-	pendix A.7.3.	258
214	mately 13000 pages, spanning 1816-1882.	4.1.1 Data Quality and Human Verification	259
215	• <i>Quebec Gazette</i> , held by the Bibliothèque	In our benchmark, we use Open LLM Genimi to ex-	260
216	et Archives Nationales du Québec, was pub-	tract the datasets from PDF files, which inevitably	261
217	lished in Québec City, and is bilingual in	causes hallucination issues and fabricates details	262
218	French-English: 51 newspaper issues, approxi-	extrinsic to the source text. To guarantee the reli-	263
219	mately 200 pages, spanning 1765-1807.	ability of our benchmark and mitigate the risk of	264
220	• <i>Revista Economica</i> , held by the Biblioteca	hallucination inherent in synthetic data generation,	265
221	Virtual de Prensa Histórica of the Ministry of	we use a rigorous <i>human-in-the-loop</i> quality as-	266
	Culture of Spain, was published in Havana	surance protocol. Specifically, we subjected the	267
		generated CQA triples to a manual verification pro-	268
		cess. Human annotators reviewed the samples with	269

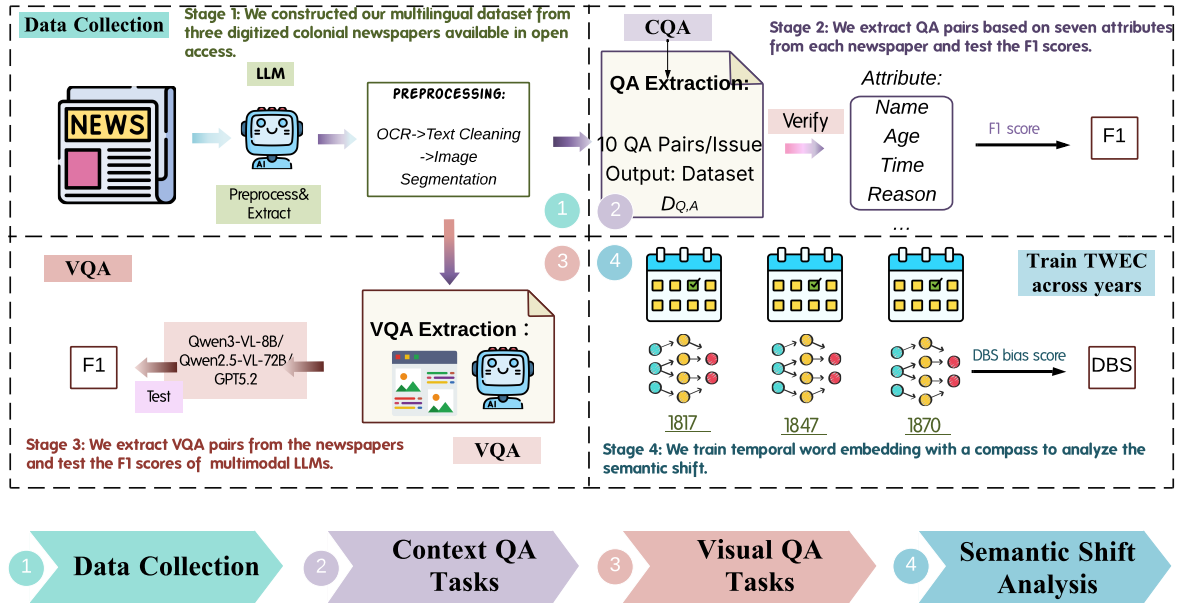


Figure 2: We conduct data collection using LLM (stage 1) and then study on multilingual event extraction (CQA tasks in stage 2 and VQA tasks in stage 3) and semantic shift tasks (stage 4).

a specific criterion: the answer a must be strictly derivable solely from the evidence provided within the context (C), without recourse to external information. This validation step effectively filtered out ungrounded responses and hallucinated facts, ensuring that our benchmark evaluates a model’s capability to reason over specific context from the newspaper.

4.2 Results

Table 1 reports the performance of different models on multilingual question answering tasks in Dutch, French, and Spanish. The scores are averaged across seven question types. Overall, larger models perform better than smaller ones in all languages. The 1.5B models show limited performance, especially on reasoning-related categories. Performance improves substantially when scaling to 72B parameters. Across all three languages, GPT-4o achieves the highest average score. It consistently outperforms open-weight models and Claude 3.5. This indicates strong multilingual and cross-task generalization. For Dutch, performance increases steadily with model size. GPT-4o obtains the best average score, followed by Claude 3.5 and Qwen2-72B. Smaller models lag behind, particularly on Reason and Position questions. For French, all models perform better than in Dutch. GPT-4o again achieves the highest average score. Large open models, such as Qwen2.5-72B, perform

competitively but remain below proprietary models. For Spanish, overall performance is the highest among the three languages. GPT-4o achieves the best result, while Qwen2-72B and Claude 3.5 show similar strong performance. Small models again perform significantly worse.

In summary, the results show a clear effect of model scale and training quality. Large models perform better across languages and question types, with GPT-4o consistently achieving the strongest results.

4.3 Visual Question Answering Tasks

In addition to CQA tasks, we test whether LLMs can answer questions when the input is a newspaper image. In this setting, the model uses the image as context and produces an answer directly. To support this, we introduce a visual question answering (VQA) benchmark that captures the multimodal nature of colonial newspapers and reduces reliance on noisy text digitization. By asking multimodal LLMs to answer questions from raw document images, we test whether they can connect language with visual text and recover meaning without using error-prone OCR systems. The statistics of the dataset are presented in Table 8. We evaluate the current state-of-the-art multimodal LLMs, such as Qwen2.5-VL-72B and closed-form LLM GPT5.2.

Model	Name	Age	Time	Reason	Position	Reward	Other	AVG
Dutch								
Qwen2-1.5B	0.616	0.416	0.497	0.372	0.322	0.550	0.472	0.464
Llama 3-8B	0.601	0.392	0.487	0.353	0.311	0.502	0.452	0.443
Qwen2-72B	0.745	<u>0.812</u>	0.672	0.512	0.683	0.612	0.632	0.667
Claude 3.5	<u>0.800</u>	<u>0.812</u>	<u>0.788</u>	<u>0.572</u>	<u>0.721</u>	<u>0.700</u>	<u>0.673</u>	<u>0.724</u>
GPT-4o	0.819	0.836	0.792	0.580	0.741	0.716	0.685	0.738
English-French								
Qwen2-1.5B	0.723	0.752	0.654	0.512	0.519	0.831	0.590	0.654
Llama 3-8B	0.712	0.723	0.632	0.502	0.520	0.801	0.581	0.639
Qwen2.5-72B	0.821	0.912	0.872	<u>0.621</u>	<u>0.517</u>	<u>0.877</u>	0.623	0.749
Claude 3.5	<u>0.824</u>	<u>0.934</u>	0.901	0.601	0.506	0.887	<u>0.643</u>	<u>0.756</u>
GPT-4o	0.831	0.935	<u>0.899</u>	0.642	0.527	0.887	0.653	0.768
Spanish								
Qwen2-1.5B	0.653	0.624	0.570	0.373	0.478	0.677	0.596	0.567
Llama 3-8B	0.642	0.614	0.564	0.321	0.473	0.664	0.573	0.550
Qwen2-72B	0.890	0.921	0.900	<u>0.710</u>	<u>0.701</u>	0.831	<u>0.712</u>	<u>0.809</u>
Claude 3.5	<u>0.900</u>	<u>0.927</u>	0.922	0.670	0.682	<u>0.832</u>	0.702	0.805
GPT-4o	0.901	0.931	<u>0.921</u>	0.715	0.707	0.857	0.728	0.823

Table 1: Evaluation results on contextual question answering tasks are reported using F1 scores. The best and second-best results are highlighted in bold and underline, respectively. Dataset details are provided in Table 7.

4.3.1 Results

Table 2 reports the results on visual question answering tasks across Dutch, English–French, and Spanish newspapers. Scores are averaged over all available question types. Overall performance is low for all models, especially for Dutch and English–French. Small multimodal models such as Qwen2.5-VL-7B show very limited ability to answer questions from historical newspaper images. Larger models perform better but still struggle in several categories. GPT-5.2 achieves the highest average score in all three language settings. It consistently outperforms open multimodal models across most question types. Large open models such as Qwen2.5-VL-72B and LLaVA-OV-72B perform competitively but remain behind GPT-5.2. Spanish yields much higher scores than Dutch and English–French. This trend is consistent across all models and suggests that visual and textual cues are easier to extract in Spanish newspapers. In summary, VQA on historical newspapers is challenging. Larger models perform much better, with GPT-5.2 achieving the best results across all languages.

5 Semantic Shift analysis

To capture the semantic shift in the colonial newspaper, we go beyond static event extraction and perform a diachronic analysis. The period from the eighteenth to the nineteenth century includes major social and political changes. For example, the abolition of slavery in the Dutch colonies in 1863 provides a key case for studying language change over time in low-resource colonial archives.

5.1 Measures

To mathematically rigorously quantify the semantic shift in racializing language, we introduce the *diachronic bias shift* (DBS) metric, which draws upon the Word Embedding Association Test (WEAT) framework (Caliskan et al., 2017). By projecting target identity terms (e.g., "freedman") onto an axis defined by opposing attribute sets, DBS explicitly tracks the transition of the racialized subject from an object of capital to an object of social anxiety.

Specifically, the DBS metric measures the temporal evolution of the semantic distance between a specific target group and opposing attribute concepts. We define two distinct attribute sets $\{A, B\}$

Model	Name	Age	Time	Reason	Feature	Other	AVG
Dutch							
Qwen2.5-VL-7B	0.083	0.142	0.093	-	0.060	0.168	0.109
Qwen3-VL-8B	0.142	0.312	0.135	-	0.151	0.200	0.188
LLaVA-OV-72B	0.143	0.374	0.143	-	0.142	<u>0.210</u>	0.202
Qwen2.5-VL-72B	<u>0.153</u>	<u>0.478</u>	<u>0.163</u>	-	<u>0.152</u>	0.203	<u>0.230</u>
GPT-5.2	0.240	0.676	0.273	-	0.429	0.283	0.380
English-French							
Qwen2.5-VL-7B	0.102	0.013	0.145	0.092	0.132	0.204	0.115
Qwen3-VL-8B	0.212	<u>0.043</u>	0.232	0.113	0.232	0.298	0.188
LLaVA-OV-72B	0.223	0.032	0.210	<u>0.133</u>	0.252	<u>0.332</u>	0.197
Qwen2.5-VL-72B	<u>0.234</u>	<u>0.043</u>	<u>0.243</u>	0.132	<u>0.263</u>	0.324	<u>0.207</u>
GPT-5.2	0.284	0.083	0.481	0.209	0.275	0.384	0.286
Spanish							
Qwen2.5-VL-7B	0.342	-	0.463	0.362	0.132	0.323	0.324
Qwen3-VL-8B	0.645	-	0.693	0.501	0.241	0.392	0.494
LLaVA-OV-72B	0.682	-	0.732	0.492	0.242	0.403	0.510
Qwen2.5-VL-72B	<u>0.672</u>	-	<u>0.743</u>	0.521	<u>0.282</u>	<u>0.452</u>	<u>0.534</u>
GPT-5.2	0.710	-	0.776	<u>0.500</u>	0.306	0.471	0.553

Table 2: Evaluation on Visual Question Answer tasks, "-" indicates that there is no dataset extracted from the newspaper. The best and second-best results are highlighted in bold and underline, respectively.

representing opposing semantic framings. For example, we set archaic vocabulary of coerced servitude ($\mathcal{A}_{servitude}$) from the emerging modern framework of free market contractuality (\mathcal{B}_{market}).

- $\mathcal{A}_{servitude} = \{slave, toil, forced, master\}$
- $\mathcal{B}_{market} = \{wage, contract, strike, pay\}$

Let $\mathbf{w}_t \in \mathbb{R}^d$ be the aligned vector representation of a target term w at time step t . The bias score $Bias_t(w)$ is defined as the difference between the mean cosine similarity of the target word to \mathcal{A} and \mathcal{B} .

$$Bias_t(w) = \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \cos(\mathbf{w}_t, \mathbf{a}_t) - \frac{1}{|\mathcal{B}|} \sum_{b \in \mathcal{B}} \cos(\mathbf{w}_t, \mathbf{b}_t) \quad (1)$$

where \mathbf{a}_t and \mathbf{b}_t denote the vector representations of the attribute words at time t . The sign of $Bias_t(w)$ reveals the dominant semantic framing in the corpus for a given era. To construct temporally aligned vector spaces for DBS analysis, we employ the Temporal Word Embeddings with a Compass (TWEC) framework (Di Carlo et al., 2019) and generate a temporally aligned vector

space for each annual data slice from 1816 to 1882 (5.2).

5.2 Temporal Word Embeddings with a Compass (TWEC)

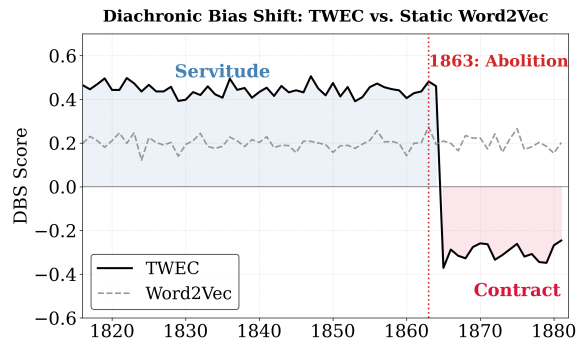


Figure 3: DBS score for "labor". The trajectory reveals a semantic phase transition. The vertical dashed line marks the 1863 Abolition of Slavery, where the discursive framing shifts from "Servitude" to "Contract".

TWEC operates in two stages: (1) training independent word embedding models for each temporal slice, and (2) aligning these embeddings into a shared coordinate system using stable "compass" words.

Let $\mathcal{D} = \{D_{t_1}, D_{t_2}, \dots, D_{t_n}\}$ denote our diachronic corpus partitioned into n temporal slices,

where D_{t_i} represents the document collection from time period t_i . For each slice D_{t_i} , we train an independent Word2Vec model M_{t_i} using the SkipGram architecture with negative sampling (Algorithm 1).

Algorithm 1 TWEC Training Framework

Require: Diachronic corpus $\mathcal{D} = \{D_{t_1}, \dots, D_{t_n}\}$, embedding dimension d , window size w , minimum frequency f_{min}

Ensure: Aligned embeddings in shared coordinate system: $\{\mathbf{W}^{(t_1)}, \dots, \mathbf{W}^{(t_n)}\}$

- 1: **Stage 1: Train Independent Slice Models**
- 2: **for** each time slice $t_i \in \{t_1, \dots, t_n\}$ **do**
- 3: $D_{t_i}^{clean} \leftarrow \text{Tokenize}(D_{t_i})$
- 4: $\mathcal{V}_{t_i} \leftarrow \{w : \text{count}(w, D_{t_i}) \geq f_{min}\}$
- 5: $\mathbf{W}_0^{(t_i)} \leftarrow \text{SkipGram}(D_{t_i}^{clean}, \mathcal{V}_{t_i})$
- 6: $\mathcal{M} \leftarrow \mathcal{M} \cup \{(t_i, M_{t_i}, \mathbf{W}_0^{(t_i)})\}$
- 7: **end for**
- 8: **Stage 2: Procrustes Alignment**
- 9: Select reference time: $t_{ref} \leftarrow t_1$
- 10: $\mathbf{W}^{(t_{ref})} \leftarrow \mathbf{W}_0^{(t_{ref})}$ {Reference unchanged}
- 11: **for** each slice $t_i \in \{t_2, \dots, t_n\}$ **do**
- 12: Extract compass vectors from reference:
- 13: $\mathbf{C}_{ref} \leftarrow [\mathbf{w}_1^{(t_{ref})}, \dots, \mathbf{w}_{|C|}^{(t_{ref})}] \in \mathbb{R}^{|C| \times d}$
- 14: Extract compass vectors from current slice:
- 15: $\mathbf{C}_{t_i} \leftarrow [\mathbf{w}_1^{(t_i)}, \dots, \mathbf{w}_{|C|}^{(t_i)}] \in \mathbb{R}^{|C| \times d}$
- 16: Compute alignment matrix via Procrustes (Algorithm 2):
- 17: $\mathbf{Q}_{t_i} \leftarrow \text{Procrustes}(\mathbf{C}_{t_i}, \mathbf{C}_{ref})$
- 18: Align all embeddings: $\mathbf{W}^{(t_i)} \leftarrow \mathbf{W}_0^{(t_i)} \cdot \mathbf{Q}_{t_i}$
- 19: **end for**
- 20: **return** $\{\mathbf{W}^{(t_1)}, \dots, \mathbf{W}^{(t_n)}\}$

In the second stage, we align the embeddings in stage 1 into a shared coordinate system using stable "compass" words (refer to Algo 2).

5.3 Results

Figure 4 illustrates the temporal evolution of bias scores for target terms. The commodification set \mathcal{A}_{prop} is $\{sale, price, buy, owner\}$ and the criminalization set \mathcal{B}_{threat} is $\{riot, dangerous, vagrant, steal\}$ (refer to the Appendix A.4). We observe significant shifts in the semantic framing of several terms across the 19th century. (1) **Slave terms** (*slaaf*) exhibited strong commodification bias (Bias $> +0.3$) in 1817–1840, reflecting economic objectification. Post-1847, we observe a gradual shift toward neutral framing (Bias ≈ 0), corresponding with

abolitionist discourse. (2) **Free persons** (*vrijman*) maintained consistently positive commodification bias throughout the period, though decreasing from +0.42 (1817) to +0.28 (1870), suggesting evolving conceptions of free labor. (3) **Indigenous terms** (*inboorling*): Demonstrated oscillation between framing categories, with peak criminalization bias (−0.31) in 1847 during labor shortage periods, reverting to near-neutral by 1870. (4) **Workers** (*arbeider*): Showed increasing commodification bias over time (+0.15 to +0.37), reflecting the rise of wage labor discourse post-emancipation.

The results reveal that racialization is a fluid process of linguistic recalibration, where the decline of formal enslavement did not eliminate bias but instead transferred commodification and criminalization frameworks onto "workers" to sustain colonial power structures.

5.4 Discussion

To ensure robust results, we conduct several analyses. First, we evaluate the alignment quality of TWEC (Sec.5.4.1). Next, we compare TWEC with a static word2vec model (Sec.5.4.2).

5.4.1 Alignment Quality

We test the alignment error during 1817-1840, 1817-1847, 1817-1870. Table 3 shows the Procrustes alignment quality for each temporal slice relative to the 1817 reference.

Metric	1817→1840	1817→1847	1817→1870	Mean
Alignment error	0.087	0.112	0.134	0.111
Compass words used	487	473	456	472

Table 3: Procrustes Alignment Quality (Mean Euclidean Distance)

5.4.2 TWEC vs Word2vec

To test the quality of TWEC vs static Word2vec, we compare the TWEC and static word embedding training. As shown in Figure 3, the term "arbeid" (Labor in Dutch) exhibits a distinct structural break, transitioning sharply from the positive domain associated with Servitude ($\mathcal{A}_{servitude}$) to the negative domain associated with Contract (\mathcal{B}_{market}) immediately following the 1863 threshold. This "crossover" trajectory quantitatively confirms that the abolition of slavery precipitated a rapid discursive realignment, effectively decoupling labor from the lexicon of coerced bondage and integrating it into the framework of free market economics. In contrast, Word2Vec serves as a stability check.

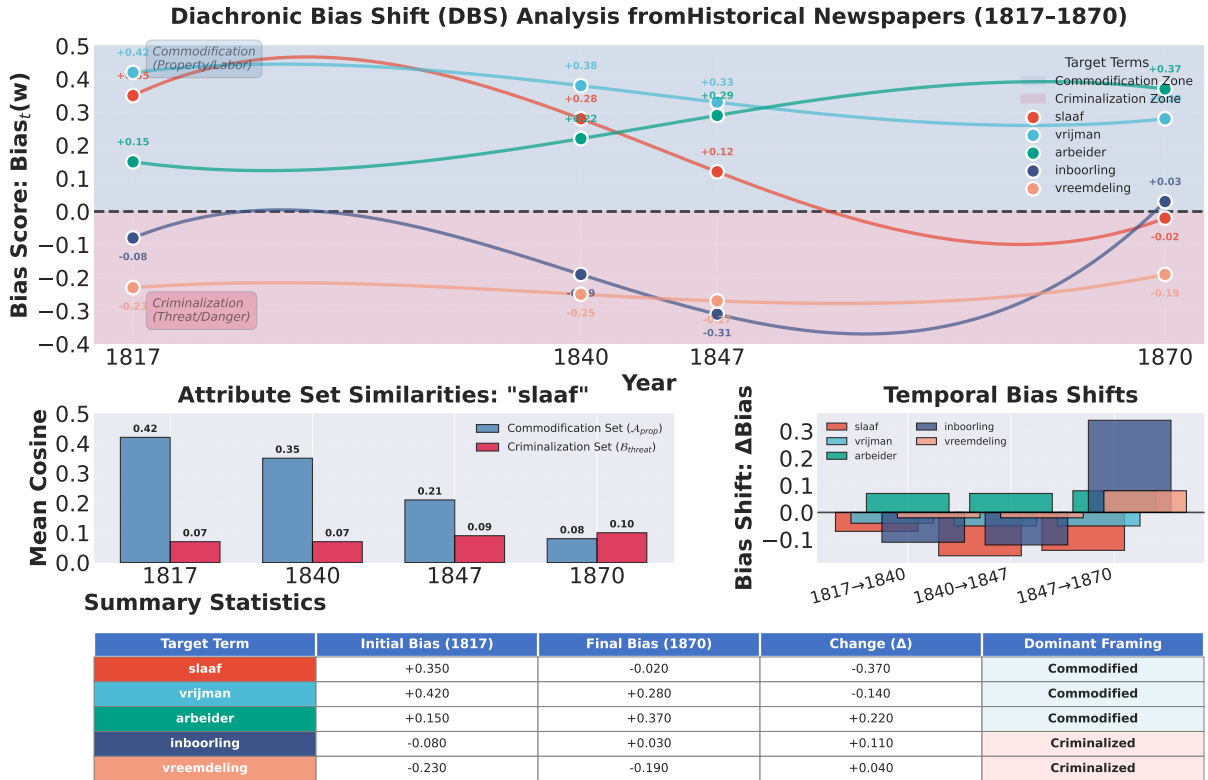


Figure 4: Diachronic Bias Shift trajectories for target demographic terms in Dutch newspapers. Positive values indicate commodification bias (association with property/labor), while negative values indicate criminalization bias (association with threat/danger) (In this case, $\mathcal{A}_{prop} = \{sale, price, buy, owner\}$, $\mathcal{B}_{threat} = \{riot, dangerous, vagrant, steal\}$). Shaded regions denote semantic framing zones.

This demonstrates the advantages of TWEC in temporal semantic analysis.

5.4.3 Statistical Significance

Target Term	1817→1840	1840→1847	1847→1870
slaaf	< 0.001***	0.023*	< 0.001***
vrijman	0.089	0.156	0.041*
arbeider	0.012*	0.003**	< 0.001***
inboorling	0.034*	< 0.001***	0.067
vreemdeling	0.234	0.412	0.378

Table 4: Statistical significance of bias shifts (p -values indicate the probability of observing the measured bias shift by chance alone. Values below 0.05 represent statistically significant temporal semantic change.). * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$.

We assess the statistical significance of observed bias shifts using bootstrap resampling ($N = 1000$ iterations). Table 4 presents p -values for temporal comparisons.

The term 'slaaf' exhibited highly significant semantic shifts across all periods ($p < 0.001$ for 1817→1840 and 1847→1870), providing strong evidence for diachronic change in the seman-

tic framing of enslaved persons. In contrast, 'vreemdeling' showed no significant shifts (all $p > 0.05$), suggesting stable criminalization framing throughout the corpus period.

6 Conclusion

In this paper, we study event extraction and semantic shift in colonial newspapers. We find that LLMs perform well on contextual question answering tasks. However, their accuracy drops sharply when the context is images of original newspapers. This shows that current multimodal models still struggle with archaic language, even though they perform well on modern text. Our semantic shift analysis also shows that racialization changed after the end of formal slavery. Bias moved from enslaved people to "workers" and indigenous terms through the new language of control and criminalization. This suggests that many modern language biases remain as a form of "semantic debt" from these historical power structures.

7 Limitations

We study racializing language in several languages from the eighteenth and nineteenth centuries. One limitation of our work is that we focus only on Dutch, English, French, and Spanish newspapers. Many other colonial languages and regions are not included. In future work, we plan to use more historical newspapers and study a wider range of languages. In addition, our analysis depends on large language models. These models may reflect biases from their modern training data, which can influence the results. Finally, our semantic shift analysis is limited by data size. Some years contain fewer documents, which may reduce the reliability of the observed language changes. Future work could use larger and more balanced datasets to address this issue.

8 Ethical Considerations

Studying texts about the history of slavery raises important ethical concerns (Rickford, 2016). The enslaved people described in the newspapers used in this study lived centuries ago, so issues of personal privacy and data protection do not directly apply. However, the newspapers contain many examples of racist and demeaning language, and this language can be harmful and distressing, even when presented in a historical context. Preserving this linguistic characteristic was necessary for us to assess how accurately LLMs process and analyze these historical documents. Nonetheless, this requires appropriate interpretation and careful handling.

Acknowledgments

References

H. Samy Alim. 2016. Introducing raciolinguistics: Racializing language and languaging race in hyper-racial times. In *Raciolinguistics*, pages 1–30.

Nadav Borenstein, Natália da Silva Perez, and Isabelle Augenstein. 2023a. Multilingual event extraction from historical newspaper adverts. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10304–10325.

Nadav Borenstein, Karolina Stańczak, Thea Rolskov, Natacha Klein Käfer, Natália da Silva Perez, and Isabelle Augenstein. 2023b. Measuring intersectional biases in historical documents. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2711–2730.

Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.

Ruirui Chen, Chengwei Qin, Weifeng Jiang, and Dongkyu Choi. 2024. Is a large language model a good annotator for event extraction? In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 17772–17780.

Komala Subramanyam Cherukuri, Pranav Abishai Moses, Aisa Sakata, Jiangping Chen, and Haihua Chen. 2025. Large language models for oral history understanding with text classification and sentiment analysis. *arXiv preprint arXiv:2508.06729*.

James Costa and Zorana Sokolovska. 2025. Historical foundations: Some threads for integrating and interrogating historiography in critical sociolinguistics. *Critical sociolinguistics: dialogues, dissonances, developments*.

Natália da Silva Perez and Nadav Borenstein. 2025. Annotating "privacy" to train semi-supervised event extraction models for historical newspapers. *Current Research in Digital History*, 8.

John Dagdelen, Alexander Dunn, Sanghoon Lee, Nicholas Walker, Andrew S Rosen, Gerbrand Ceder, Kristin A Persson, and Anubhav Jain. 2024. Structured information extraction from scientific text with large language models. *Nature communications*, 15(1):1418.

Valerio Di Carlo, Federico Bianchi, and Matteo Palmonari. 2019. Training temporal word embeddings with a compass. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 6326–6334.

Xinya Du, Alexander M Rush, and Claire Cardie. 2021. Grit: Generative role-filler transformers for document-level event entity extraction. In *Proceedings of the 16th conference of the european chapter of the association for computational linguistics: main Volume*, pages 634–644.

Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644.

Jack Grieve, Sara Bartl, Matteo Fuoli, Jason Grafmiller, Weihang Huang, Alejandro Jawerbaum, Akira Murakami, Marcus Perlman, Dana Roemling, and Bodo Winter. 2025. The sociolinguistic foundations of language modeling. *Frontiers in Artificial Intelligence*, 7:1472411.

Igor Grossmann, Matthew Feinberg, Dawn C Parker, Nicholas A Christakis, Philip E Tetlock, and William A Cunningham. 2023. Ai and the transformation of social science research. *Science*, 380(6650):1108–1109.

603	Kunal Handa, Alex Tamkin, Miles McCain, Saffron	Fayuan Li, Weihua Peng, Yuguang Chen, Quan Wang,	657
604	Huang, Esin Durmus, Sarah Heck, Jared Mueller,	Lu Pan, Yajuan Lyu, and Yong Zhu. 2020. Event	658
605	Jerry Hong, Stuart Ritchie, Tim Belonax, and 1 oth-	extraction as multi-turn question answering. In <i>Find-</i>	659
606	ers. 2025. Which economic tasks are performed with	<i>ings of the Association for Computational Linguistics:</i>	660
607	ai? evidence from millions of claude conversations.	<i>EMNLP 2020</i> , pages 829–838.	661
608	<i>arXiv preprint arXiv:2503.04761</i> .		
609	Jakob Hauser, Daniel Kondor, Jenny Reddish, Majid	Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. 2020.	662
610	Benam, Enrico Cioni, Federica Villa, James Bennett,	A joint neural model for information extraction with	663
611	Daniel Hoyer, Pieter Francois, Peter Turchin, and 1	global features. In <i>Proceedings of the 58th annual</i>	664
612	others. 2024. Large language models’ expert-level	<i>meeting of the association for computational linguis-</i>	665
613	global history knowledge benchmark (hist-llm). <i>Ad-</i>	<i>tics</i> , pages 7999–8009.	666
614	<i>vances in Neural Information Processing Systems</i> ,	Mingyu Derek Ma, Xiaoxuan Wang, Po-Nien Kung,	667
615	37:32336–32369.	P Jeffrey Brantingham, Nanyun Peng, and Wei Wang.	668
616	Mark J Hill and Simon Hengchen. 2019. Quantifying	2024. Star: boosting low-resource information ex-	669
617	the impact of dirty ocr on historical text analysis:	traction by structure-to-text data generation with	670
618	Eighteenth century collections online as a case study.	large language models. In <i>Proceedings of the AAAI</i>	671
619	<i>Digital Scholarship in the Humanities</i> , 34(4):825–	<i>Conference on Artificial Intelligence</i> , volume 38,	672
620	843.	pages 18751–18759.	673
621	Frederik Hogenboom, Flavius Frasinca, Uzay Kaymak,	Susanna Mäkinen. 2022. Stability and variation in the	674
622	and Franciska De Jong. 2011. An overview of event	genre of runaway slave notices in american new-	675
623	extraction from text. <i>DeRiVE@ ISWC</i> , pages 48–57.	papers 1704–1865. <i>Ennen ja nyt: Historian ti-</i>	676
624	Mark Humphries, Lianne C Leddy, Quinn Downton,	<i>etosanomat</i> , 22(3):83–87.	677
625	Meredith Legace, John McConnell, Isabella Murray,	Simon P Newman. 2022. <i>Freedom seekers: Escaping</i>	678
626	and Elizabeth Spence. 2024. Unlocking the archives:	<i>from slavery in restoration London</i> . University of	679
627	large language models achieve state-of-the-art perfor-	London Press.	680
628	mance on the transcription of handwritten historical	Thien Nguyen and Ralph Grishman. 2018. Graph con-	681
629	documents. <i>Available at SSRN</i> .	volutional networks with argument-aware pooling for	682
630	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam	event detection. In <i>Proceedings of the AAAI Confer-</i>	683
631	Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,	<i>ence on Artificial Intelligence</i> , volume 32.	684
632	Akila Welihinda, Alan Hayes, Alec Radford, and 1	Andrew Piper and Sophie Wu. 2025. Evaluating large	685
633	others. 2024. Gpt-4o system card. <i>arXiv preprint</i>	language models for narrative topic labeling. In <i>Pro-</i>	686
634	<i>arXiv:2410.21276</i> .	<i>ceedings of the 5th International Conference on Nat-</i>	687
635	Tristan Karch, Jakhongir Saydaliev, Isabella Di Lenardo,	<i>ural Language Processing for Digital Humanities</i> ,	688
636	and Frederic Kaplan. 2025. Llm agents for interac-	pages 281–291.	689
637	tive exploration of historical cadastre data: frame-	John R Rickford. 2016. <i>Raciolinguistics: How lan-</i>	690
638	work and application to venice. <i>Computational Hu-</i>	<i>guage shapes our ideas about race</i> . Oxford Univer-	691
639	<i>manities Research</i> , 1:e11.	sity Press.	692
640	Andrey Kutuzov, Lilja Øvrelid, Terrence Szymanski,	Christophe Rigaud, Antoine Doucet, Mickaël Coustaty,	693
641	and Erik Velldal. 2018. Diachronic word embed-	and Jean-Philippe Moreux. 2019. Icdar 2019 com-	694
642	dings and semantic shifts: a survey. <i>arXiv preprint</i>	petition on post-ocr text correction. In <i>2019 interna-</i>	695
643	<i>arXiv:1806.03537</i> .	<i>tional conference on document analysis and recogni-</i>	696
644	Viet Dac Lai, Minh Van Nguyen, Heidi Kaufman, and	<i>tion (ICDAR)</i> , pages 1588–1593. IEEE.	697
645	Thien Huu Nguyen. 2021. Event extraction from	Natália da Silva Perez. 2025. Racialised language in	698
646	historical texts: A new dataset for black rebellions.	colonial newspaper advertisements during the eigh-	699
647	In <i>Findings of the Association for Computational</i>	teenth and nineteenth centuries. In <i>The Routledge</i>	700
648	<i>Linguistics: ACL-IJCNLP 2021</i> , pages 2390–2400.	<i>Handbook of Information History</i> , pages 95–109.	701
649	Maria Levchenko. 2025. Evaluating llms for histori-	Arno Simons, Michael Zichert, and Adrian Wüthrich.	702
650	cal document ocr: A methodological framework for	2025. Large language models for history, philosophy,	703
651	digital humanities. <i>arXiv preprint arXiv:2510.06743</i> .	and sociology of science: Interpretive uses, method-	704
652	Simon Levis Sullam, Giorgia Minello, Rocco Tripodi,	ological challenges, and critical perspectives. <i>arXiv</i>	705
653	and Massimo Warglien. 2022. Representation of jews	<i>preprint arXiv:2506.12242</i> .	706
654	and anti-jewish bias in 19th century french public	Matthew Sims, Jong Ho Park, and David Bamman. 2019.	707
655	discourse: Distant and close reading. <i>Frontiers in big</i>	Literary event detection. In <i>Proceedings of the 57th</i>	708
656	<i>Data</i> , 4:723043.	<i>annual meeting of the association for computational</i>	709
		<i>linguistics</i> , pages 3623–3634.	710

Rachele Sprugnoli and Sara Tonelli. 2019. Novel event detection and classification for historical texts. *Computational Linguistics*, 45(2):229–265.

Qwen Team and 1 others. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2(3).

Konstantin Todorov and Giovanni Colavizza. 2022. An assessment of the impact of ocr noise on language models. *arXiv preprint arXiv:2202.00470*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Crina Tudor, Beáta Megyesi, and Robert Östling. 2025. Prompting the past: Exploring zero-shot learning for named entity recognition in historical texts using prompt-answering llms. In *Proceedings of the 9th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature (LaTeCH-CLfL 2025)*, pages 216–226.

Jianshu Weng and Bu-Sung Lee. 2011. Event detection in twitter. In *Proceedings of the international aaai conference on web and social media*, volume 5, pages 401–408.

Wei Xiang and Bang Wang. 2019. A survey of event extraction from text. *IEEE Access*, 7:173111–173137.

James Zhang, Wouter Haverals, Mary Naydan, and Brian W Kernighan. 2024a. Post-ocr correction with openai’s gpt models on challenging english prosody texts. In *Proceedings of the ACM Symposium on Document Engineering 2024*, pages 1–4.

Shibingfeng Zhang and Giovanni Colavizza. 2025. Named entity recognition of historical text via large language model. *arXiv preprint arXiv:2508.18090*.

Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Pan, and Lidong Bing. 2024b. Sentiment analysis in the era of large language models: A reality check. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3881–3906.

A Appendix

A.1 Algorithm

We present the alignment algorithm used in our Algo 1.

Algorithm 2 Orthogonal Procrustes Alignment

Require: Source matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$, target matrix $\mathbf{B} \in \mathbb{R}^{m \times d}$

Ensure: Optimal orthogonal matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$

1: Center matrices:

$$2: \quad \boldsymbol{\mu}_A \leftarrow \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i, \quad \boldsymbol{\mu}_B \leftarrow \frac{1}{m} \sum_{i=1}^m \mathbf{b}_i$$

$$3: \quad \mathbf{A}^c \leftarrow \mathbf{A} - \mathbf{1}_m \boldsymbol{\mu}_A^T, \quad \mathbf{B}^c \leftarrow \mathbf{B} - \mathbf{1}_m \boldsymbol{\mu}_B^T$$

4: Compute cross-covariance: $\mathbf{M} \leftarrow (\mathbf{B}^c)^T \mathbf{A}^c$

5: Singular Value Decomposition: $\mathbf{M} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$

6: Optimal rotation: $\mathbf{Q} \leftarrow \mathbf{U} \mathbf{V}^T$

7: **return** \mathbf{Q}

A.2 TWEC Model Configuration

Table 5 presents the hyperparameters used for TWEC training. For Word2vec training, we use the gensim² to train the word2vec separately for each year. All the training is conducted on an A100 80G GPU.

Table 5: TWEC Hyperparameters

Parameter	Value
Embedding dimension (d)	100
Context window (w)	5
Minimum word frequency (f_{\min})	3
Skip-gram negative samples	5
Training epochs (T)	15
Learning rate	0.025
Compass vocabulary size ($ \mathcal{C} $)	500
CV stability threshold ($\tau_{\text{stability}}$)	0.5

A.3 Computational Complexity

The TWEC framework has time complexity $\mathcal{O}(n \cdot T \cdot |V| \cdot d)$ for Stage 1, where n is the number of time slices, T is the number of training epochs, $|V|$ is the vocabulary size, and d is the embedding dimension. Stage 2 Procrustes alignment has complexity $\mathcal{O}((n-1) \cdot d^3)$ due to SVD computation. DBS computation for each target word requires $\mathcal{O}(n \cdot (|\mathcal{A}| + |\mathcal{B}|) \cdot d)$ operations.

A.4 DBS attribute sets

Table 6 presents the attribute sets we used in the computation of DBS scores.

A.5 Prompt of Data Generation

We present the prompt template we used to generate the CQA and VQA pairs from the newspaper.

²<https://pypi.org/project/gensim/>

Attribute	Wordlist
$\mathcal{A}_{servitude}$	{ <i>slave, toil, forced, master</i> }
\mathcal{B}_{market}	{ <i>wage, contract, strike, pay</i> }
\mathcal{A}_{prop}	{ <i>sale, price, buy, owner</i> }
\mathcal{B}_{threat}	{ <i>riot, dangerous, vagrant, steal</i> }
\mathcal{A}_{nature}	{ <i>water, flow, stream, bank</i> }
$\mathcal{B}_{politics}$	{ <i>vote, law, war, rights</i> }

Table 6: Words used for DBS evaluation.

Context-based Question Answer Pairs Prompt Template

User Prompt:

You are a historian specializing in 19th-century history. I will provide you with an original runaway slave advertisement (Context) in newspapers in different languages. Based on this text, please extract key information and convert it into ten context-based question answer (CQA) pairs from every newspaper.

Task requirements:

1. Identify key attributes mentioned in the advertisement (e.g., name, age, time, reason, position, reward, etc.).
2. For each attribute, generate a question in its original language.
3. First prioritize slavery-related content, including runaway enslaved persons, reward notices, sales, ownership, and physical descriptions. Only if insufficient, use other colonial advertisements or notices.
4. Extract the corresponding text span from the original advertisement as the answer.
5. Answers must be short and factual. Do not infer, summarize, or use external knowledge.
6. The output format must be a Python tuple list in the form: (Context, Question, Answer).

Visual-based Question Answer Pairs Prompt Template

User Prompt:

You are a historian specializing in 19th-century history. I will provide you with an original runaway slave advertisement (Context) in newspapers in different languages. Based on this text, please analyze the given newspaper page image and extract some high-quality Visual-Question-Answer (VQA) pairs.

Task requirements:

1. Identify key attributes mentioned in the advertisement (e.g., name, age, time, reason, position, reward, etc.).
2. For each attribute, generate a question in its original language.
3. First prioritize slavery-related content, including runaway enslaved persons, reward notices, sales, ownership, and physical descriptions. Only if insufficient, use other colonial advertisements or notices.
4. Extract the corresponding text span from the original advertisement as the answer.
5. Answers must be short and factual. Do not infer, summarize, or use external knowledge.
6. Page_Num must be an integer (starting from 1).
7. The output format must be a Python tuple list in the form: (Page_Num, Context, Question, Answer).

A.6 Prompt of Data Evaluation

Data Evaluation Prompt Template

User Prompt:

You are a historian specializing in 19th-century history. Answer the question using only information explicitly stated in the provided context.

Task requirements:

1. The answer must be in the same language as the question.
2. Do not translate between languages.
3. If the answer appears verbatim in the context, copy it exactly.
4. If multiple possible spans exist, choose the shortest correct span.
5. Do not explain, paraphrase, or add extra words.
6. The answer should be as short and factual as possible.

A.7 Dataset setting

A.7.1 Data Preprocessing for VQA

Building upon the filtered corpus, we further construct a visual question answering (VQA) dataset to incorporate visual and layout information. For this purpose, newspaper pages are randomly sampled, and VQA pairs are generated by aligning questions and answers with the corresponding full-page newspaper images. Unlike the text-only QA setting, the contextual input in the VQA dataset consists of the entire newspaper page containing the original textual evidence, rather than a localized text excerpt. This design aims to better reflect realistic historical newspaper reading scenarios. In total, the final dataset comprises 195 VQA instances, which are used in all experiments reported in this paper.

A.7.2 Dataset Processing for CQA

Our dataset is derived from digitized *De Curaçaosche Courant* newspapers published between 1816 and 1882. *Quebec Gazette* newspapers published between 1765 and 1807. *Revista Economica* published between 1878 and 1882. We employ Gemini 2.5 to automatically extract question-answer (QA) pairs from the newspaper text (The prompt is presented in Sec ??). In the ini-

Attribute	Dutch	French	Spanish
Name	605	208	163
Age	158	31	14
Time	163	70	58
Reason	21	4	15
Position	94	76	64
Reward	36	8	3
Other	250	104	169

Table 7: Data statistics for contextual question answer tasks.

Attribute	Dutch	French	Spanish
Name	86	38	32
Age	19	4	-
Time	24	43	20
Reason	-	5	11
Feature	9	4	7
Other	63	38	96

Table 8: Data statistics for visual question answer tasks.

tial sampling stage, we adopt a year-wise strategy: for each year in the target period, one newspaper is randomly selected, from which ten QA pairs are extracted, and the resulting QA pairs predominantly concern slavery and closely related topics.

We observe a substantial decline in slavery-related content in newspapers published after 1863, the year in which slavery was officially abolished in the Netherlands. As a consequence, newspapers from 1864 to 1882 yield very few relevant QA instances. To maintain topical coherence and data consistency, we therefore restrict our dataset to the period 1816–1863, excluding later years from subsequent analysis.

A.7.3 Comparison Models

To evaluate the proposed benchmark across varying scales of compute and accessibility, we selected a diverse set of baseline models categorized into three distinct tiers. First, to test the efficacy of lightweight, locally deployable models, we include Qwen2-1.5B (Team et al., 2024), llama 3-8B (Touvron et al., 2023). Second, to assess the capabilities of high-performance open-weight architectures, we evaluate Qwen2-72B, analyzing whether open-sourced large language models can approximate the sociological reasoning required for this task. Finally, we establish a performance upper bound using state-of-the-art proprietary models, including

806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833

834 Claude 3.5 (Handa et al., 2025), GPT-4o (Hurst
835 et al., 2024), thereby enabling a critical compari-
836 son between closed-source reasoning engines and
837 accessible alternatives in the context of historical
838 text analysis.

839 For VQA tasks, we include these compari-
840 son models to provide a fair and representative
841 benchmark across model scale, architecture, and
842 training paradigms in multilingual VQA. Specif-
843 ically, Qwen2.5-VL-7B and Qwen3-VL-8B rep-
844 resent strong small-to-mid-scale open-source vi-
845 sion-language models, allowing us to assess
846 performance under realistic resource constraints.
847 Qwen2.5-VL-72B and LLaVA-OV-72B serve as
848 large-scale open-source baselines, enabling anal-
849 ysis of how scaling impacts multilingual and
850 reasoning-heavy visual tasks. GPT-5.2 is included
851 as a state-of-the-art proprietary model, providing
852 an upper-bound reference for current multimodal
853 capabilities. Together, this selection spans diverse
854 model sizes, training recipes, and openness levels,
855 ensuring that observed performance differences re-
856 flect fundamental model capabilities rather than
857 idiosyncratic design choices, and allowing us to
858 contextualize gains across languages and question
859 types.

860 A.8 Case Study

861 To provide qualitative evidence complementary to
862 the quantitative evaluation, we present one cor-
863 rectly answered example for both context-based
864 question answering (CQA) and visual question an-
865 swering (VQA), selected directly from our dataset.

Contextual QA example

Context:

“de Negerin Markita slavin van den Heer
John Harrison”

Question: Wie is de eigenaar van de slavin
Markita?

Answer: De Heer John Harrison

867 In this text-only setting, the model accurately
868 retrieves an explicitly stated factual detail from
869 a short and unambiguous context. The question
870 requires direct fact lookup without multi-sentence
871 reasoning, resulting in an exact match between the
872 predicted answer and the ground truth.

Visual QA example

Context:

[Image:] 1851-08-30-Decura...p4.jpg

Question: 26 jaren

Answer: 26 jaren

873 For the VQA case, the model successfully
874 grounds its answer in the visual content of the doc-
875 ument. The age information is clearly stated in the
876 scanned newspaper page, enabling accurate visual
877 localization and answer extraction despite noise
878 commonly present in historical document images.
879

880 A.9 Corpus Statistics

881 Our diachronic corpus consists of historical news-
882 paper articles from *Dutch* spanning 1816–1882,
883 partitioned into temporal slices as shown in Table 9,
884 Table 10.

Table 9: Corpus Statistics by Temporal Slice (1816-1850)

Year	Documents	Tokens	Vocabulary
1816	38	235770	10201
1817	53	308366	12442
1818	52	303416	12036
1819	52	358698	13522
1820	53	376178	14464
1821	52	380717	14337
1822	52	388701	15120
1823	37	286532	11841
1824	30	236484	10229
1825	51	392029	15409
1826	52	392634	15325
1827	48	351255	14044
1828	52	357560	13744
1829	52	335191	13180
1830	52	357683	13590
1831	52	368451	12611
1832	50	338282	12325
1833	51	328237	11432
1834	52	282555	10785
1835	38	200114	8261
1836	53	258490	10063
1837	50	247870	9711
1839	50	241960	9368
1840	52	265107	10010
1841	49	250195	9355
1842	52	259500	9586
1843	46	231405	9042
1844	51	269157	9635
1845	48	248169	10007
1846	52	262630	9781
1847	51	266436	9734
1848	49	269460	9926
1849	51	318787	10981
1850	49	301943	11406

Table 10: Corpus Statistics by Temporal Slice (1851-1882)

Year	Documents	Tokens	Vocabulary
1851	42	264557	10156
1852	52	320806	11752
1853	53	320537	11631
1854	52	329665	12160
1855	52	320978	11995
1856	53	303065	11344
1857	52	311291	11461
1858	52	317592	11719
1859	50	313851	11691
1861	46	289846	10919
1862	52	330886	12020
1863	53	329853	12122
1864	51	310423	12107
1865	49	306312	11449
1866	50	311487	11820
1867	52	318379	11568
1868	53	312775	11610
1869	51	310094	12003
1870	54	303236	11827
1871	51	297879	11841
1872	50	283720	11295
1873	52	280279	11543
1874	48	259878	10898
1875	50	274819	10964
1876	52	294767	11304
1877	48	330662	12496
1878	50	338294	12629
1879	51	338617	12430
1880	52	339583	12202
1881	52	339275	12048
1882	52	385273	12297