

Detecting Latin in Historical Books with Varied Layout: A Multimodal Benchmark

Anonymous ACL submission

Abstract

This paper presents a novel task of extracting Latin fragments from mixed-language historical documents with varied layouts. We benchmark and evaluate the performance of large foundation models against a multimodal dataset of 753 annotated pages. The results demonstrate that reliable Latin detection with contemporary LLMs is achievable. Our study provides the first comprehensive analysis of these models’ capabilities and limitations for this task¹.

1 Introduction

Accurate language identification at a granular level within historical documents is a key component to the study of the early modern period at scale. Latin, as the primary written language of Western Europe for more than a millennium, has a unique position, gradually ceding to vernaculars at varying paces across regions and genres (Marjanen et al., 2025). Throughout this transition, Latin fragments frequently appeared within predominantly vernacular texts, in quotations, specialist terminology, and instances of code-switching. Automated extraction of diverse Latin uses in context from historical corpora is crucial to studying language evolution, the interplay between classical and modern thought, and the dissemination of ideas (Sprugnoli et al., 2024; Gorovaia et al., 2024; Perrone et al., 2021; Burns et al., 2021). However, this task poses challenges due to wide variations in Latin usage, scripts, complex page layouts, and inconsistent print and scan quality in historical book databases.

This study focuses on detecting instances of Latin within the *Eighteenth Century Collection Online* (ECCO) (Tolonen et al., 2022) corpus using both book page images and the corresponding OCR text. The lack of an existing dataset specifically

designed for multimodal and code-mixed Latin detection motivated us to create an annotated dataset for this purpose. Our dataset contains 753 pages sampled from historical documents, validated by specialists in 18th-century publishing culture to represent diverse use cases. While we focus on Latin due to its aforementioned importance for historical study, our fully benchmarked, manually annotated scenario provides a solid template for extending the method to other languages as well.

Given the complex nature of the task and factors ranging from OCR noise to varied print layouts, we explore the capabilities of modern Large Language Models (LLMs), including multimodal models (MLLMs) for this task. The way these models handle contextual information, recognize patterns in noisy data, and integrate textual cues with visual layout information has been found to help in disambiguating text in historical documents (Luo et al., 2024; Boros et al., 2024; Kanerva et al., 2025; Xie et al., 2025).

Our investigation exploits the new dataset to test a number of state-of-the-art models and approaches, and finds that reliable Latin detection in such challenging historical material is achievable. The benchmarking of different model architectures provides insight into their strengths and weaknesses when faced with the complexities inherent in the data. This work establishes a strong baseline for a novel NLP task and at the same time highlights the need for modality-aware approaches and robust evaluation frameworks in historical text analysis.

The main contributions of this article are:

- Defining the task of Latin detection in historical documents, an understudied multimodal case of language detection.
- Creating an expert-annotated dataset from 18th-century books, capturing diverse and challenging examples of Latin usage.

¹The dataset is in the supplementary materials. Both the dataset and code will be published upon acceptance.

- Developing an evaluation framework for this task, considering challenges from both textual and visual modalities.
- Conducting a systematic comparison of contemporary LLMs to assess their efficacy for this task.
- Delivering a practical Latin detection pipeline, demonstrating its readiness for downstream applications in historical research.

2 Problem Definition

We define the task of *Latin language detection in historical documents* as a two-stage classification and extraction problem, where the input consists of a scanned page image and/or its OCR transcription. The task is to automatically detect whether any segments in the text are written in Latin, and if so, to extract text of those specific segments.

Formally, given a document page D , let I_D denote its image and T_D denote its OCR'd text. A system must perform the following two subtasks:

- **Task 1 (Page-level Latin Detection):** Predict a binary label $y_D \in \{0, 1\}$, where $y_D = 1$ indicates that the page contains at least one segment in Latin, and $y_D = 0$ otherwise.
- **Task 2 (Latin Segment Extraction):** If $y_D = 1$, extract a list of text spans $S_D = [s_1, s_2, \dots, s_n]$, where each $s_i \in T_D$ is a contiguous Latin segment string.

The reason for specifying the task in two stages is mainly due to evaluation. As our final aim is to evaluate how well Latin is identified at a high level of granularity across different layouts, we measure task 2 in terms of per-page precision and recall. However, this type of measurement does not adequately cater to cases where Latin is not present on a page. Thus, task 1 has been designed to give us this more general information on e.g., whether the models are excessive in detecting Latin where there is none. Requiring segment-level output as strings rather than image regions aligns better with tasks most MLLMs are trained on, and enables simpler performance comparisons between the input modalities. More details on the metrics used will be given in section 5.

3 Related work

Latin in NLP Given its historical importance, Latin has attracted considerable attention within

the NLP community (e.g., Sprugnoli et al., 2024; Schulz and Keller, 2016; Sprugnoli et al., 2022; Gorovaia et al., 2024; Perrone et al., 2021; Burns et al., 2021), though much of this research has centered on small, clean corpora of ancient literary texts. While some recent studies have ventured into Early Modern mixed-language documents (Stüssi and Ströbel, 2024; Volk et al., 2024), these also predominantly rely on manually curated and annotated data. In contrast, our work focuses on the foundational task of Latin *discovery*: detecting Latin within extensive, unedited, and noisy digitized collections like ECCO (Tolonen et al., 2022). This computational approach aims to detect Latin in a vast corpora, while the identified fragments can subsequently be analyzed using a range of established NLP tools developed for classical languages (Johnson et al., 2021; Burns, 2023; Straka and Straková, 2020; Kupari et al., 2024).

Code-mixed Language Detection From a methodology perspective, identifying Latin segments within British publications is a code-mixed language detection task (Aguilar et al., 2020). While extensive research in this area has focused on contemporary informal texts (Barman et al., 2014; Zhang et al., 2018), its application to historical documents, with challenges like archaic syntax, lexicon, and spelling, has been less explored (Schulz and Keller, 2016; Volk et al., 2022). Detecting classical languages in these complex historical contexts has traditionally involved rule-based systems and supervised machine learning approaches, notably Conditional Random Fields (CRFs) (Schulz and Keller, 2016; Sterner and Teufel, 2023; Volk et al., 2022). Alongside these, robust statistical tools like Lingua (Stahl, 2021) offer effective general language identification with support for mixed language. Given the recognized potential of modern LLMs to navigate linguistic nuances and noisy data, our work investigates their capacity to enhance detection performance.

LLMs for Historical Documents Recent LLMs, particularly Multimodal variants (MLLMs), have shown considerable potential in historical document analysis, demonstrating top performance in tasks like OCR, named entity recognition, and general document understanding from historical sources (Bai et al., 2025; Luo et al., 2024; Boros et al., 2024; Kanerva et al., 2025; Backer and Hyman, 2025; Xie et al., 2025), and in assessing general historical knowledge (Hauser et al., 2024). De-

spite these advancements, a significant gap persists for more specialized, complex applications. Specifically, there is a notable lack of dedicated benchmarks and systematic exploration for the fine-grained, page-level multimodal detection and extraction of embedded secondary languages (e.g., Latin) (Aguilar et al., 2020; Guzmán et al., 2017). This task is demanding due to noisy scans from historical archives, diachronic language context, and orthographic variation (Volk et al., 2022). Our work contributes to this underexplored area by introducing a systematic evaluation methodology designed to be scalable also to other languages.

4 Dataset

4.1 Sampling and Annotations

Our approach to dataset construction began with a targeted sampling strategy to identify pages with a high likelihood of containing Latin text. We queried the Reception Reader database (Rosson et al., 2023), which indexed text reuses across the ECCO corpus using noise-resistant detection methods. From this, we randomly selected 800 reuse instances where one book was cataloged as Latin and the other as non-Latin. To ensure broad representation and reflect the diversity of the ECCO collection (approximately 200,000 books), each sampled page was drawn from a different book, covering varied publication dates and topics. However, ECCO’s language metadata is book-level, meaning that “Latin” books often contain significant non-Latin text like English introductions. Also, the reuse offsets only mark the textual overlap without specifying the language of the text segment.

These pages were then manually annotated. Annotators were tasked with drawing bounding boxes around all Latin fragments on the page images (see example in Figure 1). These visual annotations could later be reliably mapped to text offsets (locations within a string) using ECCO’s OCR positional data for ground truth text extraction. Our annotation guidelines stated marking all instances of Latin text semantically used as Latin. This included single Latin words if presented with explanations in a dictionary, as well as Latin found in headlines, editorial annotations, or footnotes. Conversely, to ensure clarity and consistency, in-line Roman/Latin names (of people, places, plants) and jargon were not specifically annotated as Latin but were treated as part of the surrounding language, as well as abbreviations, such as ‘etc’.

The annotation environment was Label Studio (Tkachenko et al., 2020-2025). The primary annotation was performed by three scholars familiar with Latin. Following this, an expert in historical texts meticulously reviewed and validated all annotations to ensure accuracy and consistency.

4.2 Dataset Characteristics

In total, 753 pages were annotated, with 623 identified as containing Latin. An expected finding during annotation was the frequent presence of other languages, such as French, German, and Greek, highlighting the dataset’s challenging multilingual nature beyond simple Latin-English code-mixing.

To contextualize model performance and to better understand the dataset’s composition, we divided the sample pages into language integration categories. Each category represents a specific way in which Latin is used in 18th-century British books and how it relates to English-language text. Depending on their content, pages containing Latin text were assigned to one or multiple categories, with some frequently appearing together (e.g., bilingual editions and footnotes), while others are mostly exclusive (e.g., full Latin). By categorizing our evaluation dataset in this way, we can assess the performance of our Latin detection approach within different contexts of language integration. The following list defines the categories used:

1. Full Latin: Pages entirely written in Latin, without any other language present.
2. Bilingual Editions (Latin / English): Pages that contain both the original Latin text and its English translation.
3. Direct quotations in Latin: Pages where Latin phrases or sentences are quoted verbatim within an otherwise English text.
4. Independent Latin sections: Pages with original Latin text by the author, accompanied by English text on the same page.
5. Code Switching: Pages where the text alternates between Latin and English within the span, often for stylistic or rhetorical purposes.
6. Dictionaries: Pages with entries that define individual Latin words, often with translations or explanations in another language.
7. Footnotes: Pages with annotations or footnotes that provide explanations of Latin words or phrases used in the main text.

Table 1 shows the frequency of each annotation category within our dataset. Figures 3 and 4 in

	Category	Count
1.	Full Latin	134
2.	Bilingual Editions	65
3.	Direct quotes	258
4.	Independent Latin sections	169
5.	Code-switching	109
6.	Dictionaries	19
7.	Footnotes	69

Table 1: Frequencies of annotation categories.

DISCONTENTED with his presnt condition, and defirous to be any thing but what he is, he wishes himself one of the shepherds. He then catches the idea of rural tranquillity; but soon discovers how much happier he should be in these happy regions, with **LYCORIS** at his side.

*Hic gelidi fontes, hic mollia prata, Lycori:
Hic nemus : hic ipso tecum consumerer ævo.
Nunc infans amor duri me Martis in armis;
Tela inter media, atque adversos detinet hostes.
Tu procul a patria (nec sit mihi credere) tantum
Alpinas, ab dura, nives, & frigore Rheni
Me sine sola vides. Ab te ne frigora ledant!
Ab tibi ne teneras glacies secet aspera plantas!*

Figure 1: An example of an annotated Latin fragment.

the Appendix A show category examples. This form of page-level annotation reduces the need for costly instance-level labeling, which is particularly challenging for Latin due to expertise requirements and high annotation volume. Moreover, it supports context-rich evaluation aligned with the page-level structural nature of our detection task inputs.

5 Evaluation Setting

5.1 OCR Post-Correction and Normalization

Evaluating models on the ECCO corpus is complicated by significant OCR quality discrepancies: modern models with vision capabilities often produce cleaner text than ECCO’s original OCR, while text-based models may or may not replicate the noise in their input. Such differences make direct string-based comparisons problematic and distort evaluation. To ensure meaningful assessment across all model types, we post-correct both the ground-truth Latin segments and the full input page texts. This OCR post-correction is performed using the OpenAI o1 model (Jaech et al., 2024) with a specialized prompt from (Kanerva et al., 2025).

Even after the post-correction, residual noise and other variation still remain in the data. Thus, for token-based evaluation, we apply a more traditional rule-based pre-processing pipeline to both predicted and reference strings. This deterministic pipeline, informed by our extensive experience with OCR data and domain-specific knowledge, targets common superficial textual variations. The pipeline includes Unicode normalization, ligature replacement, lowercasing, digit removal, de-hyphenation, and punctuation stripping. Subsequent to these cleaning operations, the strings are tokenized into word-level units. More details of the processing steps could be found in the Appendix B.1.

5.2 Metrics

The goal of Task 1 is to detect whether a page has Latin on it. We measure this by reporting precision, recall, and F1-score in percentage. To evaluate the Latin segment extraction performance in Task 2, we calculate precision, recall, and F1 score in percentage based on token-level matches between model predictions and ground truth. A fuzzy matching mechanism is applied to pair predicted and reference tokens one by one. A match is considered valid if the token-level edit-distance is not larger than a threshold proportion θ compared to the ground-truth token length, which serves as a tunable hyperparameter. This approach accounts for minor lexical variations and OCR-induced distortions in single tokens, offering a more flexible and robust evaluation compared with exact token matching. The pseudocode of the fuzzy matching algorithm is shown in the Appendix B.2. Overall metrics are averaged across the full evaluation set.

6 Method

Our evaluation investigates the application of general instruction-following LLMs, particularly multimodal variants, for Latin segment extraction from historical documents. We propose a unified, prompt-based pipeline designed to be both practical for real-world deployment and robust for systematically and fairly evaluating the capabilities of diverse LLMs on this task.

Unified Prompting Strategy We employ a single, high-level instructive prompt designed to elicit responses that inherently address both sub-tasks within a unified and simply formatted output. This approach simplifies interaction with the models and

the subsequent processing of their outputs, thereby contributing to the overall ease of application.

This unified prompt asks the LLM to extract all Latin segments to a list, including single Latin words, without further instruction. The distinction in our experiments lies solely in the input provided to this consistent prompt, where the specific prompts are shown in the Appendix C:

- **Text-only input:** The LLM receives the OCR-extracted and post-corrected text, appended to the prompt.
- **Image-only input:** The LLM receives the page image, with the prompt guiding it to identify Latin script based on visual signal.
- **Multimodal input:** The LLM receives both the scanned page image and the post-corrected OCR-extracted text, allowing it to leverage both visual and textual information for the task.

Structured Output and Postprocessing The LLMs are instructed to output their predictions as a list of Latin segments, which directly corresponds to the output requirement for Task 2. The presence of a non-empty list implicitly indicates the presence of Latin script on the page (Task 1, $y_D = 1$), while an empty list indicates its absence ($y_D = 0$).

Model-Agnostic Compatibility Because the method does not rely on any model-specific architecture or training, it can be directly applied to a wide range of general-purpose foundation language models. This makes the approach particularly suitable for scalable deployment across large historical corpora with variable OCR quality and image-text alignment conditions.

7 Experiments

7.1 Experiment Setup

Model Selection To explore how modern instruction-tuned language models handle the new task of Latin segment detection and extraction in noisy multimodal historical documents, we benchmark a representative suite of LLMs across modalities, scales, and architectures. The model selection follows three guiding principles: (i) in the absence of dedicated multimodal benchmarks for historical language understanding in documents, we refer to leaderboard performance on **DocVQA** (Tito et al., 2021) and comprehensive open evaluations such as **OpenCompass** (Contributors, 2023) and **MMMU**

(Yue et al., 2024); (ii) we prioritize lightweight to medium-scale models (7B-32B) to better reflect realistic research use cases in historical academic and low-resource scenarios. Specifically, the most notable selected models include (additionally, we evaluate further models detailed in the Appendix E.2):

- **Qwen2.5-VL series** (32B, 7B) (Bai et al., 2025): flagship open-source multimodal LLMs with strong document understanding and fine-grained visual grounding. We include both 32B and 7B variants, and ablate visual inputs via text-only configurations (Qwen2.5-32B/14B/7B).
- **InternVL3 series** (14B, 9B, 8B) (Zhu et al., 2025): notable academic multimodal LLMs with a two-stage visual encoder adding transformer design. We ablate the visual inputs in the 9B model via testing its pure language counterpart Internlm3-8B-Instruct.
- **Mistral-Small-3.1-24B-Instruct** (Mistral AI, 2025): an efficient model using Mixture-of-Experts (MoE). Notably, for its strong text performance and reasoning capabilities supports an extended 128k token context window.
- **DeepSeek-R1-Distill variants** (Qwen-32B, Qwen-14B, Llama-8B) (DeepSeek-AI, 2025): pure-text distilled models focused on reasoning, probing extraction capability with thinking capability but without visual cues.

Baseline We employ **Lingua** (Stahl, 2021), a statistical language identifier based on character n-gram modeling, and the only off-the-shelf tool we found that supports token-level Latin detection in mixed-language text. While not designed for noisy OCR, it offers a useful traditional baseline to contextualize the difficulty of our task and the potential advantages and drawbacks of LLM-based approaches. More configuration of the baseline is in Appendix D.

Implementation Details All LLMs are deployed on a unified inference backend using the vLLM framework (version 0.8.5) with server mode to simulate realistic client-side usage. Inference is performed on a 4×A100 (40GB) GPU node, with no data parallelism and a batch size of 1 for all cases, resulting in an average of 16 GPU hours per model. We use each model’s default generation parameters without further tuning. For the Mistral model, we specifically load the Mistral-Small-3.1-24B-Instruct-2503 version checkpoint. The edit-

distance threshold θ is set to 0.2 based on empirical evidence, discussed further in the Appendix E.3.

7.2 Main Results and Comparisons

The main experimental results are shown in Table 2. The traditional Lingua language identifier proves to be a surprisingly strong baseline. In the page-level detection task (Task 1), Lingua achieves a strong F1 score of 92.89. In the more challenging segment extraction task (Task 2), it still attains 76.74 F1 using simple token-level identification.

Only the largest instruction-tuned foundation models such as Qwen2.5-VL-32B and DeepSeek-R1-Distill-32B consistently outperform Lingua across both tasks. Qwen in particular demonstrates strong document-oriented OCR capabilities, excelling even in vision-only settings. Interestingly, the DeepSeek-R1-Distilled version of Qwen does worse than plain Qwen on Task 1, but better on Task 2. Looking in more details at the numbers, the distilled version seems to be simply optimising for recall at the cost of precision on the page-level metric, while then being more accurate in the actual segment extraction. Smaller models, like InternVL3-14B, can still improve over Lingua when vision is used to guide text modeling, but fail in vision-only configurations, illustrating that OCR robustness is not easy to achieve.

Performance tends to improve with model scale, especially within the same model family. Larger models have been shown to more effectively memorize and generalize low-resource language phenomena, consistent with neural scaling laws (Gordon et al., 2021; Kaplan et al., 2020). However, size alone does not guarantee strong performance. Some smaller models illustrate that fusion architecture and fine-tuning strategy can outweigh raw scale. Both Qwen2.5-VL-7B and InternVL3-8B share the Qwen2.5-7B language backbone, yet Qwen2.5-VL-7B’s distinct native dynamic image resolution support and instruction-tuning on multilingual document-centric corpora gain advantages. A direct comparison between Qwen2.5-7B and InternLM3-8B reveals that Qwen’s textual pretraining data composition and cross-lingual instruction tuning confer superior robustness to OCR noise and low-resource Latin patterns, beyond what can be achieved by architecture scale alone.

Across all model families, multimodal variants (I+T) consistently outperform their unimodal counterparts under the same model size. This is especially evident in the InternVL3 series, where

I-only models fall well below their I+T variants on both tasks. This highlights the value of aligning visual features with textual reasoning. However, text-only models can still perform competitively, particularly those with fine-grained OCR-style pre-training, such as Qwen2.5. Additionally, most models cannot extract Latin solely from images, where only Qwen2.5-VL-32B matches its text-only performance, further highlighting historical document OCR challenges and our dataset contributions on dealing with OCR texts and annotations.

7.3 Performance by Category

Figure 2 shows the results on the five top models considering the coverage of different input modalities, over different page categories, which were specified in Section 4.2. As can be seen from the left-hand side of the figure, for pages with only a single type, there is a large disparity between category difficulty: models yield almost perfect performance for full-Latin and bilingual pages while struggling with code-switching and to a lesser degree with dictionaries. The difference between modalities is also evident: image and combined inputs remain robust across categories, while text-only models demonstrate more variety. In particular, text-based Qwen fails in code-switch detection: it does not find any correct Latin text at all in 5 pages out of 9 in this category, showing its shortcomings in nuanced detail extraction. Since footnotes almost never exist in isolation, we are unable to directly report performance relating to them. In addition, considering multi-label cases, the performance metrics should be dominated by the performance on the longer text category. For example, the performance on independent Latin sections plus dictionaries is better than dictionaries alone. However, when looking at the performance of footnotes in conjunction with other short text types such as code-switching and quotations, we can still clearly find that the performance of footnotes is even worse than that of code-switching.

7.4 Behavior on Non-Latin pages

To gain further insight into model behavior, we analyze performance on pages devoid of Latin. For Task 1, we report the recall of the negative class. For Task 2, we report false positive token rates (Table 3). It is evident from these analyses that all models over-detect Latin as compared to our ground-truth. The best performance is demonstrated by Qwen2.5 with text-only input, which

MODEL DETAILS				PAGE-LEVEL (TASK 1)			TOKEN-LEVEL (TASK 2)		
Model	Variant	Size	Mode	F1	Precision	Recall	F1	Precision	Recall
Lingua	-	-	T	92.89	86.85	99.84	76.74	76.64	80.01
Qwen2.5	VL	32B	I+T	94.09	88.97	99.84	84.15	86.21	84.25
	VL	32B	I	94.22	89.45	99.52	78.86	81.31	78.97
	-	32B	T	97.68	97.29	98.07	80.53	84.60	79.46
	-	14B	T	91.54	88.46	94.86	70.12	74.6	72.04
	VL	7B	I+T	89.93	83.87	96.95	73.23	77.1	76.05
	VL	7B	I	92.86	88.04	98.23	71.16	74.38	73.36
InternVL3	-	7B	T	88.37	83.33	94.05	52.43	53.5	64.17
	-	14B	I+T	91.5	84.68	99.52	80.09	81.28	82.3
	-	14B	I	87.37	82.95	92.28	44.48	46.45	46.89
	-	9B	I+T	84.69	82.37	87.14	59.88	59.33	68.09
	-	9B	I	70.82	79.28	63.99	25.28	25.91	29.61
	-	8B	I+T	89.5	83.36	96.62	65.76	61.35	81.79
Internlm3	-	8B	I	86.67	83.68	89.87	48.58	45.91	59.2
	-	8B	T	82.33	83.01	81.67	47.49	52.3	51.32
Mistral-Small-3.1	-	24B	I+T	92.73	88.47	97.43	79.89	79.4	82.93
	-	24B	I	89.82	85.16	95.02	66.14	65.39	69.81
	-	24B	T	88.47	83.26	94.37	72.97	72.61	78.79
DeepSeek-R1-Distill	Qwen2.5	32B	T	94.05	89.42	99.2	81.87	84.37	82.96
	Qwen2.5	14B	T	91.69	85.12	99.36	77.98	81.62	78.82
	Llama-3.1	8B	T	89.74	84.01	96.3	66.59	73.06	68.38

Table 2: Experimental results on selected LLMs, compared with Lingua baseline. “I” indicates image input, “T” indicates text input, and “I + T” refers to the combination of both modalities.

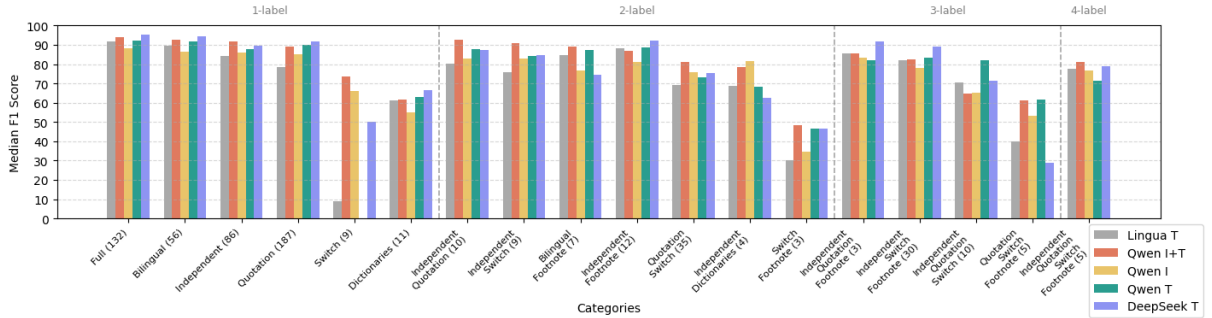


Figure 2: Result on different category labels for 5 top-performing models: median of the token-level F1-score for each page. The number of instances for each label is shown in parentheses. We filter out the subsets with fewer than 3 cases to ensure the statistical validity of the subset results.

correctly refrains from outputting Latin on 87% of the non-Latin pages, with a low proportion of Latin token output on misclassified pages. This could suggest that it benefits from a stronger sensitivity to linguistic context, enabling more cautious predictions in ambiguous cases.

The statistical baseline Lingua does remarkably worse in this evaluation, misidentifying Latin on over 70% of non-Latin pages. This is primarily be-

cause the ngram-based model struggles with deeper linguistic nuances and context, and outputs false positive tokens on most pages. Since our goal is to process a collection of 200K books, a large number of false positives would pose a significant problem.

In general, the number of Latin tokens identified on non-Latin pages remains small, pointing to the possibility of using further filtering or thresholding approaches in downstream tasks.

MODEL DETAILS				PAGE	TOKEN
Model	Variant	Size	Mode	Neg. Recall	FP Rate
Lingua	-	-	T	28.24	2.65
Qwen2.5	VL	32B	I+T	41.22	3.41
	VL	32B	I	44.27	1.94
	-	32B	T	87.02	0.53
DeepSeek	Qwen2.5	32B	T	44.27	4.84

Table 3: Analysis on non-Latin pages: Negative Class Recall (in pages) and False Positive Rate (percentage of tokens falsely identified as Latin on the whole page).

7.5 Qualitative Evaluation and Error Analysis

Our qualitative evaluation and error analysis are informed by the ECCO dataset’s characteristics and ambiguities in defining Latin. These factors complicated annotation, influenced model predictions, and set practical limits on achievable evaluation scores. For examples and example-oriented descriptions of the issues outlined here, see Appendix E.1.

First, poor image quality in ECCO, along with complex page layouts, such as multi-column text, footnotes, marginalia, and varied fonts or spacing (see Figures 3 and 4, Appendix A), often led to inaccurate OCR and fragmented transcriptions, hampering reliable annotation, adding noise to even the OCR-post-corrected page texts.

The error analysis of Latin-containing segments shows lower performance in categories like code-switching, dictionaries, and footnotes. This is caused by both the brevity of Latin snippets in these contexts and the fact that they are more likely to be affected by severe OCR errors, the latter often due to complex layouts and challenging font types or sizes on poor-quality scans. Figure 6 in Appendix E.1 shows an example of how OCR errors can lead to discrepancies between our ground truth and predictions.

Content-wise, an interesting definitional mismatch is observed, both decreasing precision as well as increasing false positive rates. In our ground-truth annotation guidelines, we specified that Roman named entities, common anglicized Latin phrases (e.g. “e.g.”, “etc.”, etc.) and jargon within otherwise English text were not to be considered Latin. However, these are frequently identified by the models as Latin (for an example, see Figure 7 in Appendix E.1). In terms of performance, these discrepancies often account for the full difference between ground truth and prediction, meaning that with a slight change in definitions, the model performance would be even stronger.

PROMPT CONFIG		PAGE			TOKEN		
Non-Latin	Abbrev	F	P	R	F	P	R
		94.09	88.97	99.84	84.15	86.21	84.25
✓		93.37	87.82	99.68	84.28	85.34	85.42
✓	✓	93.01	87.31	99.52	84.26	85.74	84.89

Table 4: Prompt modification experiments.

7.6 Prompt Modification Experiments

To assess the impact of prompt fine-tuning, we conduct a study on the best model (Qwen2.5-VL-32B) by (i) explicitly instructing it to output an empty list when no Latin is detected (Non-Latin column) and (ii) adding explicit instructions to not extract common Latin abbreviations found in other languages (Abbrev column). As shown in Table 4, these modifications yield only marginal changes, indicating small gains in some cases compensated by losses in others. This is particularly interesting concerning the abbreviations, as it means that the model has a strong innate definition of what it considers Latin, which cannot be easily overcome by prompting.

8 Conclusion

This paper presented a novel task and a dataset for Latin extraction from early-modern book pages. We systematically evaluated diverse foundational models and found that this task can be solved with excellent performance, without fine-tuning. However, such performance can be achieved only with bigger models (32B parameters).

Our results show that 94% F1 performance on page-level detection can be achieved with only image input. This has practical implications since a visual model can be used even in collections that have no OCR, or whose OCR is not of sufficient quality. In this case, OCR or post-OCR correction can be performed only on pages preselected by a visual model, which would save computation cost.

We also found that semantic analysis is crucial to distinguish actual Latin phrases from usages of Latin words in English text. Thus we argue for the usage of large foundation models, or a pipeline of visual and textual models, rather than pure statistical methods.

Further steps will include processing the whole ECCO collection and publishing a complete dataset of Latin fragments in the ECCO books. We also plan to expand our results to other Early-modern collections, such as the French BNF collection.

Limitations

All our work is performed on a single collection. Results from this single corpus (18th-century British ECCO texts) should be interpreted with caution, as potential corpus-specific biases may affect broader generalizability, highlighting the importance of future cross-corpus studies. Since Latin has been widely used across Europe, it would be interesting to validate our methods on other collections, e.g., the main texts written in Romance languages, such as French.

A drawback for per-category evaluation in the current dataset is the lack of per-segment type annotations, which leads to a shortage of instance-level performance analysis. These will be implemented in the next version of the dataset.

The current work does not utilize a separate validation set, a decision necessitated by the considerable difficulty and cost of annotating the historical data. Consequently, model selection was performed without an independent dataset for optimization before final testing, instead referring to the results from previous work. Incorporating a validation set is an important consideration for future extensions of this research.

All reported results are from single experimental runs in each model’s default configuration. This approach was a necessary compromise due to the significant computational costs and GPU resource constraints associated with evaluating the diverse range of large-scale models investigated. Future work could explore the impact of running variability by conducting multiple trials and tuning seeds, temperatures, and other hyperparameters.

It would also be interesting to compare the off-the-shelf LLMs with some models trained or fine-tuned specifically for the Latin extraction task. However, given the surprising performance of general-purpose LLMs on this task, the practical usefulness of such experiments is questionable.

Ethics Statement

The underlying literary works from which our dataset is derived, sourced from 18th-century texts within the Eighteenth Century Collections Online (ECCO), are in the public domain. The compilation and sharing of our dataset, which comprises annotated excerpts and portions of page images from this collection, are conducted for research purposes under the permissions granted. We are committed to ensuring that the creation and dissemination of this dataset adhere to relevant copyright

considerations and ethical guidelines.

We used ChatGPT and Gemini for grammar and spell-checking and stylistic polishing of the draft of this manuscript. All suggestions were critically reviewed and edited by the authors to ensure factual accuracy and originality.

References

- Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. 2024. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*.
- Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, et al. 2024. Pixtral 12b. *arXiv preprint arXiv:2410.07073*.
- Gustavo Aguilar, Sudipta Kar, and Tamar Solorio. 2020. [LinCE: A centralized benchmark for linguistic code-switching evaluation](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1803–1813, Marseille, France. European Language Resources Association.
- Samuel Backer and Louis Hyman. 2025. Bootstrapping ai: Interdisciplinary approaches to assessing ocr quality in english-language historical documents. In *Proceedings of the 5th International Conference on Natural Language Processing for Digital Humanities*, pages 251–256.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Utsab Barman, Amitava Das, Joachim Wagner, and Jennifer Foster. 2014. [Code mixing: A challenge for language identification in the language of social media](#). In *Proceedings of the First Workshop on Computational Approaches to Code Switching*, pages 13–23, Doha, Qatar. Association for Computational Linguistics.
- Emanuela Boros, Maud Ehrmann, Matteo Romanello, Sven Najem-Meyer, and Frédéric Kaplan. 2024. [Post-correction of historical text transcripts with large language models: An exploratory study](#). In *Proceedings of the 8th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature (LaTeCH-CLfL 2024)*, pages 133–159, St. Julians, Malta. Association for Computational Linguistics.

750	Patrick J Burns. 2023. Latency: Synthetic	lms in historical documents: No free lunches. <i>arXiv</i>	806
751	trained pipelines for latin nlp. <i>arXiv preprint</i>	<i>preprint arXiv:2502.01205</i> .	807
752	<i>arXiv:2305.04365</i> .		
753	Patrick J Burns, James A Brofos, Kyle Li, Pramit Chaud-	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B	808
754	huri, and Joseph P Dexter. 2021. Profiling of inter-	Brown, Benjamin Chess, Rewon Child, Scott Gray,	809
755	textuality in latin literature using word embeddings.	Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.	810
756	In <i>Proceedings of the 2021 Conference of the North</i>	Scaling laws for neural language models. <i>arXiv</i>	811
757	<i>American Chapter of the Association for Computa-</i>	<i>preprint arXiv:2001.08361</i> .	812
758	<i>tional Linguistics: Human Language Technologies</i> ,		
759	pages 4900–4907.	Hanna-Mari Kristiina Kupari, Erik Henriksson,	813
		Veronika Laippala, and Jenna Kanerva. 2024. Im-	814
760	OpenCompass Contributors. 2023. Opencompass:	proving Latin dependency parsing by combining tree-	815
761	A universal evaluation platform for foundation	banks and predictions . In <i>Proceedings of the 4th</i>	816
762	models. https://github.com/open-compass/	<i>International Conference on Natural Language Pro-</i>	817
763	opencompass .	<i>cessing for Digital Humanities</i> , pages 216–228, Mi-	818
		ami, USA. Association for Computational Linguis-	819
764	DeepSeek-AI. 2025. Deepseek-r1: Incentivizing rea-	tics.	820
765	soning capability in llms via reinforcement learning.		
766	<i>Preprint</i> , arXiv:2501.12948.	Chuwei Luo, Yufan Shen, Zhaoqing Zhu, Qi Zheng, Zhi	821
		Yu, and Cong Yao. 2024. Layoutllm: Layout instruc-	822
767	Mitchell A Gordon, Kevin Duh, and Jared Kaplan. 2021.	tion tuning with large language models for document	823
768	Data and parameter scaling laws for neural machine	understanding. In <i>Proceedings of the IEEE/CVF con-</i>	824
769	translation. In <i>Proceedings of the 2021 Conference</i>	<i>ference on computer vision and pattern recognition</i> ,	825
770	<i>on Empirical Methods in Natural Language Process-</i>	pages 15630–15640.	826
771	<i>ing</i> , pages 5915–5922.		
772	Svetlana Gorovaia, Gleb Schmidt, and Ivan P.	Jani Marjanen, Tuuli Tahko, Leo Lahti, and Mikko Tolo-	827
773	Yamshchikov. 2024. Sui generis: Large language	nen. 2025. Book printing in latin and vernacular	828
774	models for authorship attribution and verification in	languages in northern europe, 1500–1800. In <i>The</i>	829
775	Latin . In <i>Proceedings of the 4th International Con-</i>	<i>Hermeneutics of Bibliographic Data and Cultural</i>	830
776	<i>ference on Natural Language Processing for Digital</i>	<i>Metadata</i> , pages 27–66. National Library of Norway.	831
777	<i>Humanities</i> , pages 398–412, Miami, USA. Associa-		
778	tion for Computational Linguistics.	Mistral AI. 2025. Mistral-small-3.1-24b-instruct-	832
		2503. https://huggingface.co/mistralai/	833
779	Gualberto A Guzmán, Joseph Ricard, Jacqueline Seri-	Mistral-Small-3.1-24B-Instruct-2503 .	834
780	gos, Barbara E Bullock, and Almeida Jacqueline	Instruction-tuned multimodal model with 24B	835
781	Toribio. 2017. Metrics for modeling code-switching	parameters, 128k context length, and vision-text	836
782	across corpora. In <i>Interspeech</i> , pages 67–71.	capabilities.	837
		Valerio Perrone, Simon Hengchen, Marco Palma,	838
783	Jakob Hauser, Daniel Kondor, Jenny Reddish, Majid	Alessandro Vatri, Jim Q Smith, and Barbara	839
784	Benam, Enrico Cioni, Federica Villa, James Bennett,	McGillivray. 2021. Lexical semantic change for an-	840
785	Daniel Hoyer, Pieter Francois, Peter Turchin, et al.	cient greek and latin. <i>Computational approaches to</i>	841
786	2024. Large language models’ expert-level global	<i>semantic change</i> , 6.	842
787	history knowledge benchmark (hist-llm). <i>Advances</i>		
788	<i>in Neural Information Processing Systems</i> , 37:32336–	David Rosson, Eetu Mäkelä, Ville Vaara, Ananth Ma-	843
789	32369.	hadevan, Yann Ryan, and Mikko Tolonen. 2023. Re-	844
		ception reader: Exploring text reuse in early modern	845
790	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard-	british publications. <i>Journal of Open Humanities</i>	846
791	son, Ahmed El-Kishky, Aiden Low, Alec Helyar,	<i>Data</i> , 9(1).	847
792	Aleksander Madry, Alex Beutel, Alex Carney, et al.		
793	2024. Openai o1 system card. <i>arXiv preprint</i>	Sarah Schulz and Mareike Keller. 2016. Code-	848
794	<i>arXiv:2412.16720</i> .	switching ubique est - language identification and	849
795	Kyle P. Johnson, Patrick J. Burns, John Stewart, Todd	part-of-speech tagging for historical mixed text . In	850
796	Cook, Clément Besnier, and William J. B. Mattingly.	<i>Proceedings of the 10th SIGHUM Workshop on Lan-</i>	851
797	2021. The Classical Language Toolkit: An NLP	<i>guage Technology for Cultural Heritage, Social Sci-</i>	852
798	framework for pre-modern languages . In <i>Proceed-</i>	<i>ences, and Humanities</i> , pages 43–51, Berlin, Ger-	853
799	<i>ings of the 59th Annual Meeting of the Association for</i>	many. Association for Computational Linguistics.	854
800	<i>Computational Linguistics and the 11th International</i>		
801	<i>Joint Conference on Natural Language Processing:</i>	Rachele Sprugnoli, Federica Iurescia, and Marco Pas-	855
802	<i>System Demonstrations</i> , pages 20–29, Online. Asso-	sarotti. 2024. Overview of the evalatin 2024 evalu-	856
803	ciation for Computational Linguistics.	ation campaign. In <i>Proceedings of the Third Work-</i>	857
		<i>shop on Language Technologies for Historical and</i>	858
804	Jenna Kanerva, Cassandra Ledins, Siiri Käpyaho, and	<i>Ancient Languages (LT4HALA)@ LREC-COLING-</i>	859
805	Filip Ginter. 2025. Ocr error post-correction with	2024, pages 190–197.	860

A Dataset Illustrations

Figures 3 and 4 show examples of Latin text categories. Both figures feature independent Latin text in the main text box at the top of the pages and footnotes at the bottom. Figure 3 is a bilingual edition of a Latin text, with an English translation directly below the Latin text at the top of the page. The footnotes in Figure 4 also include instances of code-switching and direct quotations of Latin text.

354 P. VIRG. MAR. ÆNEIDOS Lib. X.
et defertur ad antiquam domum patris Daui.
At iuventa Mezentis, mo-
nitis Jovis, ardens succedit pug-
na, troasque Teucros exan-
ta, Tyrrhenæ acies concurrunt,
atque instant viro uel, uel in-
quam omnia dâs frequentibus
bolque telis. Ille Mezentis,
velat pupis, rose prodit in
vagum æquor, cœcis f. r. ius ven-
torum, expolique ponto, per-
fert cœcis uen atque minas
velatque maribus, ipsa mœna
tonnita. Sternit Hebrum, pro-
lem Dolicaonis, huiusq; cum non
pergit Latagum f. cœcis, ne
Palmum: sed a cupit Latagum
per cœcis, neque ad resam fœ-
atque uenit, neque mœna:
fœs Palmum uenit fœs, ne-
quid fœs, neque dœs, neque fœ-
ere ejus arma, hœna, et fi-
gere cœsis, cœsis cœsis. Nec
non cœsis Phrygiæ E-
cœs, Mezentisq; æqualem
cœsis: Paridis: cum The-
ano dedit in lucem cœsis, A-
me: uenit cœsis, et Cœsis re-
gna, prœgna fœs.

Et patris antiquam Dauni defertur ad urbem:
At Jovis interea monitis Mezentius ardens
Succedit pugnae, Teucroque invadit ovantes.
Concurrunt Tyrrhenæ acies, atque omnibus
uni,
Uni odisque viro telisque frequentibus instant.
Ille, velut rupes, vastum quæ prodit in æquor,
Obvia ventorum furis, expostaque ponto,
Vim cunctam, atque minas perfert cœlicæ ma-
rrique;
Ipsa immota manens. Prolem Dolicaonis He-
brum
Sternit humi; cum quo Latagum, Palmumque
fugacem:
Sed Latagum faxo, atque ingenti fragmine
montis
Occupat os, faciemque adversam: poplite Pal-
mum
Succiso volvi segnem finit; armaque Lauro
Donat habere humeris, et vertice figere cristas.
Nec non Evantem Phrygium, Paridique Mi-
manta
Æqualem, comitemque: unâ quem nocte The-
ano
In lucem genitori Amyco dedit; et facce præ-
gnans

cutting the Deep, with prosperous Wind and Tide; and is wafted to the ancient City of his Father Daunus.

Meanwhile, by Jove's Suggestion, furious Mezentius succeeds him in the Fight, and assaults the Trojans flushed with Success. The Tuscan Troops rushed on him at once, and with all their Rage and Darts thick following each other press on him, on him alone. He stands firm as a Rock which projects into the vast Ocean, obnoxious to the Furies of the Winds, and, exposed to the Rage of the Main, endures all the Violence and Terrors of the Sky and Sea; itself unmoved remaining. He stretches on the Ground Hebrus the Son of Dolicaon, and with him Latagus and fugitive Palmus: But to Latagus with a Rock and vast Fragment of a Mountain he gives a preventing Blow on his Jaws and Face full right against him: Palmus hamstringed he suffers recreant on the Ground to roll; and gives Lausus to wear his Armour on his Shoulders, and on his Helmet's Top to fix his Plumes. Evas the Phrygian too he overthrows, and Mimas the Companion of Paris, and his Equal in Age: Whom Theano brought forth to his Father Amycus in the same Night that Queen Hecuba, the Daughter of Cisseus, preg-

N O T E S.

"Daui in lucem genitori Amyco dedit, et facce prægna fœs" i. e. observes that *creat* here is quite redundant, since the Sentence is perfect without it; be- sides

Figure 3: An example page with Latin fragments.

B Evaluation Setting Details

B.1 Preprocessing

This section details the text preprocessing pipeline for evaluation, implemented in Python, to normalize both ground truth and predicted text strings before unigram token extraction. The primary goal of this pipeline is to standardize textual representations, thereby mitigating the impact of superficial variations (e.g., from OCR noise or stylistic differences) on downstream metrics calculation. Note,

22 Q. HORATII FLACCI CARMINUM Lib. I.
Jam Cytheræ choros ducit Venus, imminente Lunâ, 5
Junctæque Nymphis Gratiæ decentes
Alterno terram quatunt pede, dum graves Cyclopus
Vulcanus ardens urit officinas.
Nunc decet aut viridi nitidum caput impedire myrto,
Aut flore, terræ quem ferunt solutæ. 10
Nunc & in umbrosis Fauno decet immolare lucis,
Seu poscat agnâ, five malit hædo.
Pallida mors æquo pulsat pede pauperum tabernas,
Regumque tures. O beate Sesti,
Vitæ summa brevis spem nos vetat inchoare longam. 15
Jam te premet nox, fabulæque Manes,
Et domus exilis Plutonia; quò simul mœaris,
Nec regna vini fortiere talis,
Nunc omnis, & mox virgines tepebunt. 20
CARMEN

5. [Jam Cytheræ choros.] The Poet here describes the Feasts of Venus, which were celebrated by young Women with Dances and Hymns in Honour of the Goddess. They began on the first of April, at the Rising of the Moon, imminente luna, and continued three Night successively. An unknown, ancient Author has thus described them:

[Jam tritum choris videri.]
Fœstas vestibus
Congregas inter cœcios
Ire per salus tuas,
Flores inter coronas,
Mortem inter calat.
Full three Nights, in joyous Vein,
Might you see the choral Train,
Hand in Hand promiscuous rove
Through thy Love-devoted Grove,
Crown'd with rosy-breathing Flowers,
Under Myrtle-woven Bowers.

6. [Gratiæ decentes.] The Graces were the most amiable Divinities of the Heathen Mythology, and the Source of all that is pleasing in Nature. The Poet calls them decentes for that Modesty and Reserve, with which they behaved themselves in these Assemblies. SAN.

[Nunc decet.] These two Verses continue the Description of the Feasts of Venus; for Flowers, and particularly Myrtle, were consecrated to that Goddess. CRUS.

that Modesty and Reserve, with which they behaved themselves in these Assemblies. SAN. The Nymphs are thus numbered by the Author already quoted:

Ruris hic erant puellæ,
Et iuvenis fontium,
Quæque sylvæ, quæque lucis,
Quæque montis incolant.
Here shall meet the blooming Maids
Of the Valleys and the Glades;
And the Nymphs, who haunt the Fountains,
And the Forests, and the Mountains. D.

7. [Gratiæ officinas.] We have here a very pretty Opposition between the Characters of Venus and Vulcan; the gay Delights of the Wife, and the laborious Employment of the Husband; who is here described working in Spring, that He might forge Thunderbolts enough for Jupiter to throw in Summer. RODELLIUS. DAC.

Figure 4: An example page with Latin fragments.

that this applies to the evaluation step only, while Latin extracting models take an input text without these steps.

For each text string, the following sequential operations are performed:

- Unicode Normalization:** Each string undergoes Unicode normalization using the normalize with "NFKD" method from Python's built-in unicodedata module. This step decomposes characters into their canonical forms, for example, separating accents from base characters, which helps in standardizing character representation.
- Ligature Replacement:** A predefined set of common ligatures is replaced with their constituent characters. Examples of replacements include 'ff' to ff, 'æ' to ae, and importantly for some historical contexts, '&' to et.
- Lowercasing:** All alphabetic characters in the string are converted to lowercase.

4. **Digit Removal:** All sequences of digits are removed from the string to avoid prediction ambiguity on digits, e.g., OCR digits in footnote notations.
5. **De-hyphenation (Word Merging):** This step addresses common OCR inconsistencies in handling end-of-line hyphens from historical documents. To ensure textual uniformity for subsequent analysis, word segments that were hyphenated, typically due to line breaks in the original source, are consistently merged into single tokens.
6. **Punctuation Stripping:** All standard punctuation marks, as defined by Python’s `string.punctuation` constant, are removed from the string.
7. **Word Tokenization:** After the above cleaning steps, each processed sequence is tokenized into a list of individual words using the `word_tokenize` function from the NLTK library (`nltk.tokenize`, version 3.9.1).

B.2 Fuzzy Matching Algorithm in Token-level Metrics

To evaluate segment correspondence, we apply a fuzzy matching algorithm to compare lists of pre-processed ground truth tokens against predicted tokens for each sample. This approach calculates Precision, Recall, and F1-score while being robust to minor textual variations. The core matching logic is outlined in Algorithm 1.

The algorithm performs a greedy, one-to-one fuzzy match: each predicted token is compared against available ground truth tokens using a match indicator function (`IsFuzzyMatch`) based on edit distance and a predefined proportion threshold θ . A match only holds when the edit distance is less than or equal to θ proportion of the length of the ground truth token string. A ground truth token can only be matched once to ensure an accurate count of distinct true positive matches. This fuzzy approach is beneficial as it offers robustness to minor textual variations that may persist even after preprocessing, leading to a more meaningful evaluation of segment correspondence.

C LLM Prompt Details

This section details the exact prompt templates employed to instruct the LLMs for the task of Latin script detection and extraction. The prompts were

adapted based on the input modality being used. In the templates below, the placeholder `{page_text}` indicates where the OCR output corresponding to the processed page image was dynamically inserted.

Image + Text

Identify and extract all segments written in Latin (e.g., Classical or Medieval Latin) from the provided image, using the accompanying OCR text as a reference. Include even single-word segments. Return the results as a list of strings in the JSON format: `[“text1”, “text2”, ...]`.

OCR Text: `{page_text}`

Image-only

Identify and extract all segments written in Latin (e.g., Classical or Medieval Latin) from the provided image. Include even single-word segments. Return the results as a list of strings in the JSON format: `[“text1”, “text2”, ...]`.

Text-only

Identify and extract all segments written in Latin (e.g., Classical or Medieval Latin) from the OCR text of an image. Include even single-word segments. Return the results as a list of strings in the JSON format: `[“text1”, “text2”, ...]`.

OCR Text: `{page_text}`

D Configuration of Baseline

For comparative language identification, we employed *Lingua* (version 2.1.0) (Stahl, 2021) as a baseline. The *LanguageDetector* was specifically configured to operate with a predefined restricted set of eight languages: English, French, German, Greek, Italian, Spanish, Portuguese, and Latin. This selection aims to encompass Latin itself and a set of the most frequently occurring languages within our target corpus ECCO (English,

Algorithm 1 Fuzzy Matching and Token Metrics Output

```
1: procedure CALCULATEFUZZYMETRICS(GT_Tokens, Pred_Tokens,  $\theta$ )
2:                                      $\triangleright$  Input: GT_Tokens, Pred_Tokens (lists of preprocessed tokens)
3:                                      $\triangleright$   $\theta$  (edit distance ratio threshold for a match)
4:                                      $\triangleright$  Output: Precision, Recall, F1-score
5:    $TP \leftarrow 0$ 
6:    $matched\_gt\_indices \leftarrow \emptyset$ 
7:   for each pred_token in Pred_Tokens do
8:     for each gt_token in GT_Tokens (with index gt_idx) do
9:       if gt_idx  $\in$  matched_gt_indices then continue
10:      end if
11:      if ISFUZZYMATCH(gt_token, pred_token,  $\theta$ ) then
12:         $TP \leftarrow TP + 1$ 
13:        Add gt_idx to matched_gt_indices
14:        break  $\triangleright$  Current pred_token matched
15:      end if
16:    end for
17:  end for
18:   $FP \leftarrow \text{length}(\text{Pred\_Tokens}) - TP$ 
19:   $FN \leftarrow \text{length}(\text{GT\_Tokens}) - TP$ 
20:   $Precision \leftarrow TP / (TP + FP)$ 
21:   $Recall \leftarrow TP / (TP + FN)$ 
22:   $F1 \leftarrow 2 \times (Precision \times Recall) / (Precision + Recall)$ 
23:  return Precision, Recall, F1
24: end procedure
```

French, German, and Greek), while also including languages present in ECCO that share orthographic or lexical similarities with Latin (Italian, Spanish, and Portuguese). Including these similar languages was intended to create a more global and robust test scenario for accurate Latin identification in ECCO.

In our pipeline, Lingua’s function to detect multiple languages within a given text (`detect_multiple_languages_of` method) was utilized on each page’s OCR output. From the resulting language segments identified by Lingua, only those substrings classified as Latin were subsequently extracted for our analysis and evaluation.

E Additional Results

E.1 Qualitative Results

More qualitative results are shown as examples to illustrate the best model’s performance and the error modes.

Figure 6 shows an example of a mismatch caused by significant OCR noise caused by poor original image quality. Here, the post-corrected OCR of our ground truth differs so much from the OCR visual or multimodal models produced during the predic-

tion process that not even our edit-distance based fuzzy ground truth matching can recover what is essentially a full match. This kind of error especially affects pages in the footnotes, code switching and dictionary categories, since the Latin texts in these categories tend to be printed in harder to detect fonts and layouts, which are additionally more likely to be affected by bad scan quality.

Figure 7 shows an example of a definitional mismatch between our annotations and the predictions. Although there is no Latin text on the page, the prediction contains the Roman names appearing in the page text.

Figure 8 shows an example of a page where the prediction contains hallucinations. The model took part of the text and translated it into Latin in the prediction, without being prompted to do so.

E.2 Other LLMs

To further contextualize the performance of contemporary foundation models on our challenging Latin discovery task, we evaluated additional models beyond those in the main comparison (Table 2). This section presents results for variants from the Pixtral and Phi-4 families. Specifically, we exam-

ined **Pixtral-12B** (Agrawal et al., 2024), an efficient Mixture-of-Experts (MoE) model, and models from the **Phi-4** family (Abdin et al., 2024), including *Phi-4-Multimodal-Instruct* (5.6B parameters) noted for its smaller footprint and a larger text-only *Phi-4* (14B parameters) variant. The results are detailed in Table 5.

The results in Table 5 reveal varied performance. Notably, some smaller multimodal models, such as *Phi-4-Multimodal-Instruct* (5.6B) particularly in image-only (I) mode, struggled significantly with the fine-grained token-level extraction task, achieving an F1 score as low as 2.48. Similarly, *Pixtral-12B* showed a substantial performance drop in its image-only configuration for token-level results. This suggests that factors like smaller parameter counts or training data less attuned to the nuances of documents and noisy OCR may limit the out-of-the-box effectiveness of certain general-purpose models for this specialized task. In contrast, the text-only *Phi-4* variant performed more competently, underscoring that model architecture with different input modalities and training focus are critical. These observations highlight the challenging nature of our proposed task.

E.3 Fuzzy Matching Threshold

The fuzzy matching threshold, θ (representing the maximum allowed normalized edit distance relative to ground truth token length), was empirically set to 0.2 for all main experiments. This choice aligns with a common heuristic of tolerating approximately “1 error in 5 characters,” suitable for OCR-derived text, and is supported by our sensitivity analysis in Figures 5. It consistently shows that while F1 scores generally increase with θ , the most substantial and steepest F1 score improvements for the majority of evaluated models are concentrated in the range leading up to $\theta \approx 0.2$, effectively compensating for common, fine-grained textual variations attributable to OCR noise. Although metrics may continue to rise beyond this point for some configurations on our dataset, we maintain $\theta = 0.2$ as a principled trade-off. A higher universal threshold could risk over-tolerating more substantial prediction errors beyond typical OCR noise, potentially prioritizing the matching of token quantity or approximate form over precise content fidelity. This could also obscure true output quality differences, especially when comparing models with varying input noise levels (e.g., image-only versus OCR-input systems).

MODEL DETAILS				PAGE-LEVEL (TASK 1)			TOKEN-LEVEL (TASK 2)		
Model	Variant	Size	Mode	F1	Precision	Recall	F1	Precision	Recall
Lingua	-	-	T	92.89	86.85	99.84	76.74	76.64	80.01
Pixtral	-	12B	I+T	89.46	84.65	94.86	69.70	71.48	74.13
			I	74.77	79.28	70.74	31.33	32.43	33.16
Phi-4	Multimodal	5.6B	I+T	76.46	85.04	69.45	39.53	42.43	46.02
	Multimodal	5.6B	I	40.84	86.91	26.69	2.48	4.13	2.78
	-	14B	T	91.64	86.24	97.75	70.14	74.17	75.12

Table 5: Experimental results for additional evaluated LLMs.

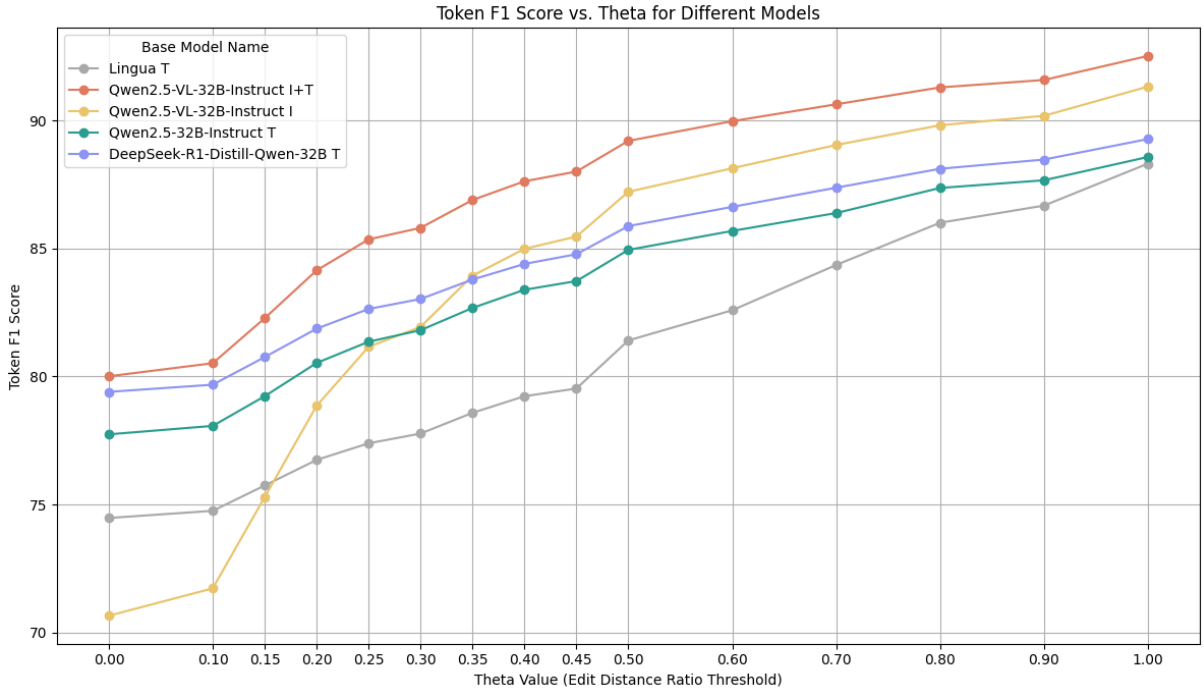


Figure 5: Token F1 scores on different θ value.

164 Of the Gods of the Heathens.

Hippotades. He dwelt in one of those seven islands which from him are called *Æolia*, and sometimes *Vulcania*. He was a skillful astronomer, and an excellent natural philosopher; he understood more particularly the nature of the winds: and because, from the clouds of smoke of the Æolian islands, he foretold winds and tempests a great while before they arose, it was generally believed that they were under his power, and that he could raise the winds or still them as he pleased. And from hence he was styled *emperor and king of the winds*, (the children of *Atræus* and *Aurora*). * *Virgil* describes,

Palæphat. de incredibil. Var. Strab. ap. Serv.
 "Nimborum in patriam loca fata turentibus Auftris,
 "Æoliam venit: Hic vasto rex Æolus antro
 "Luctantes ventos, tempestatæque sonoras
 "Imperio premit, ac vinclis & carcere frenat.
 "Illi indignantes, magno cum murmure, montis
 "Circum claustra fremunt: celsa sedet Æolus arce,
 "Sceptra tenens, mollitque animos & temperat iras.
 "Ni faciat maria, ac terras, cœlumque profundum,
 "Quippe ferant rapidi secum, verrantque per auras.
 "Sed pater omnipotens speluncis abdidit atris,
 "Hoc metuens, molemque, & montes insuper altos
 "Imposuit, regemque dedit, qui fœdere certo
 "Et premere, & laxas sciret dare jussus habenas."

Thus rag'd the Goddels, and, with fury fraught,
 The restless regions of the storms she sought:
 Where, in a spacious cave of living stone,
 The tyrant Æolus, from his airy throne,
 With pow'r imperial curbs the struggling winds,
 And founding tempests in dark prisons binds.
 This way and that, th' impatient captives tend,
 And, pressing for release, the mountains rend;
 High in the hall th' undaunted monarch stands,
 And shakes his sceptre, and their rage commands:
 Which did he not, their unresisted sway
 Would sweep the world before them in their way:
 Earth, air, and seas thro' empty space would roll,
 And heav'n would fly before the driving soul.
 In fear of this the father of the Gods
 Confin'd their fury to these dark abodes,
 And lock'd them safe, oppress'd with mountain loads;

Imposuit

GT:

palyphæ de incredibil. var. strab. ap. serv. s. r. iwhubt. urn
 in. pat. izamlocj. fâfua. urentibus. auffris. i. alolam. venit. hic.
 vato. rex. aechtes. actro. l. ludianes. vetos. tempestatæflue.
 sonoras. imperio. premit. ac. vinclis. et. carcere. frenat. illi.
 indignantes. magno. cum. murmure. montis. circum. claustra. fremunt.
 celsa. sede. jovis. arce. l. sceptra. tenens. mollitque. animos. et.
 temperat. iras. ni. faciat. maria. ac. terras. cœlumque. profundum.
 quippe. ferant. rapidi. secum. verrantque. per. auras. sed. pater.
 omnipotens. speluncis. abdidit. atris. hoc. metuens. molemque. et.
 montes. insuper. altos. lc. imposuit. regemque. dedit. qui. fœdere.
 certo. et. premere. et. laxas. sciret. dare. iussus. habenas.

Pred:

nimborum in patriam loca fata turentibus austris æoliam
 venit hic vato rex æchius antron luctantes ventos
 tempestatæque sonoras imperio premit ac vinclis et carcere
 frenat illi indignantes magno cum murmure montis circum
 claustra fremunt celsa fede t jous arcen sceptra tenens
 mollitque animos et temperat iras ni faciat maria ac terras
 cœlumque profundum quippe ferant rapidi fecum verrantque per
 auras sed pater omnipotens speluncis abdidit atris hoc
 metuens molemque et montes insuper altos impavit regemque
 dedit qui fœdere certo et premere etc laxas dare jussus
 habecas

Figure 6: An example page with Latin fragments, together with our ground truth and prediction for that page.

dent from History, that this mock Senate, this Senate in Burlesque, was compos'd of a Parcel of Scoundrels who had never seen *Pharsalia*. For can you or any one believe, that if they had been of real Senatorian Rank, *Cæsar* would have us'd them as he did, who hang'd up as many of them as fell into his Hands? But let us now see what *Sempronius* is pleas'd to reply to *Portius*.

*Not all the Pomp and Majesty of Rome
Can raise her Senate more than Cato's Presence.*

———O, my *Portius*!

*Could but I call that wond'rous Man my Father,
Would but thy Sister Marcia be propitious
To thy Friend's Vows, I might be blest in deed.*

*Port. Alas, Sempronius! would'st thou
talk of Love
To Marcia, while her Father's Life's in danger?*

Thou might'st as well court the pale trembling Vestal,

*When she beholds the Holy Flame expiring
Sempr. The more I see the Wonders of thy Race,*

The more I'm charm'd. Thou must take heed my Portius,

*The World has all its Eyes on Cato's Son.
Thy Father's Merit sets thee up to View,*

Am

GT:

Pred:

sempronius portius marcia pharsalia cato vesta

Figure 7: An example page without Latin fragments, together with our ground truth and prediction for that page.

of the Church of ENGLAND.

15

" corruption, accustoming ourselves by little and little, to
 " comprehend and bear the Majesty of God." And S.
Cyprian. " If he be made the Temple of God, I ask, Of *Epi* 73.
 " what God? If he answer, Of the Creator, he could not
 " be His Temple, because he did not believe in him. If
 " he say, Of Christ, neither can he be made His Temple,
 " because he denies Christ to be God. Or if he say, Of
 " the Holy Ghost, yet since these Three are One, how can
 " the Holy Ghost be reconciled to that Man, who is an
 " Enemy either to the Father or to the Son?"

Very and Eternal God] The most notorious Opposer of
 the Godhead of the Holy Ghost, was *Macedonius*, Patri-
 arch of *Constantinople*. The Heresy itself is called the He-
 resy of the *Pneumatomachi*, or *Fighters against the Spirit*;
 as denying the Divinity of the Holy Ghost, and asserting
 that he is only a created Energy or Power, attending upon
 and ministering unto the Son. In order to put a Stop to
 this Heresy, the first Council of *Constantinople*, to these
 Words in the *Nicene Creed*, *I believe in the Holy Ghost*,
 added, *The Lord, and Giver of Life, who proceedeth from*
the Father and the Son, Who with the Father and the Son
together is worshipped and glorified, Who spake by the Pro-
phets. See *Pearson on the Creed*, p. 325.

ARTICLE VI.

*Of the Sufficiency of the Holy Scriptures
 for Salvation.*

HOLY Scripture containeth all Things
 necessary to Salvation: So that whatso-
 ever is not read therein, nor may be proved
 thereby, is not to be required of any Man,
 that it should be believed as an Article of the
 Faith, or be thought requisite or necessary to
 Salvation. In the Name of the holy Scripture
 we do understand those Canonical Books of
 the Old and New Testament, of whose Au-
 thority was never any Doubt in the Church.

Of

GT:

Pred:

credo in spiritum sanctum dominum et vivificantem qui ex
 patre filioque procedit qui cum patre et filio adoratur et
 conglorificatur qui locutus est per prophetas

Figure 8: An example page without Latin fragments, together with our ground truth and hallucinated prediction for that page.