



HAL
open science

The Patrologia Graeca Corpus

Chahan Vidal-Gorène, Bastien Kindt

► **To cite this version:**

Chahan Vidal-Gorène, Bastien Kindt. The Patrologia Graeca Corpus. Language Resources and Evaluation Conference (LREC 2026), May 2026, Palma De Majorque, Spain. ⟨hal-05622581⟩

HAL Id: hal-05622581

<https://enc.hal.science/hal-05622581v1>

Submitted on 14 May 2026

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY-NC-SA 4.0 - Attribution - Non-commercial use - ShareAlike - International License

The Patrologia Graeca Corpus

OCR, Annotation, and Open Release of Noisy Nineteenth-Century Polytonic Greek Editions



Chahan Vidal-Gorène^{1,2}, Bastien Kindt³

¹École nationale des chartes-PSL (France) ²Calfa (France)

³UCLouvain - CIOL Institut orientaliste (Belgium)

Contact: chahan.vidal-gorene@chartes.psl.eu



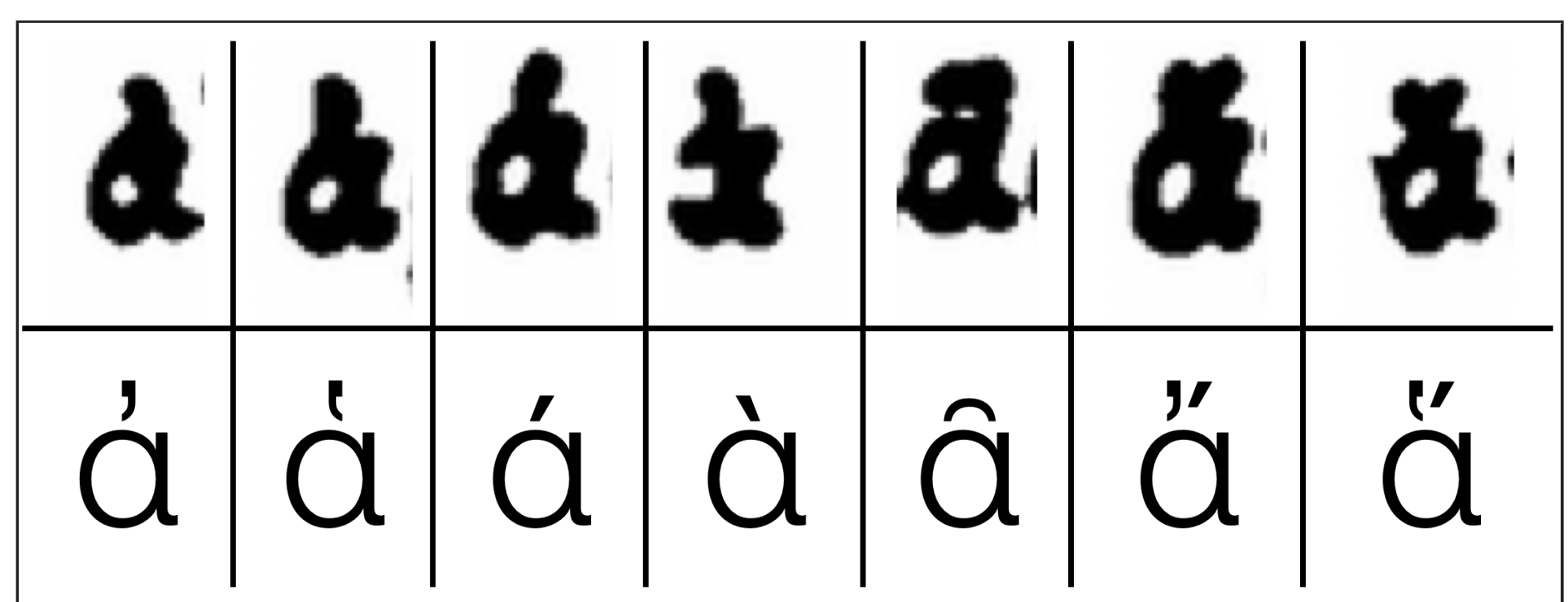
CONTEXT

Objective: OCR, lemmatization and POS-tagging of the **still-unavailable 33 volumes** of the *Patrologia Graeca* (post-classical and Byzantine texts).

Challenge: Processing noisy 19th-century editions to fill gaps left by scholarly databases like the *Thesaurus Linguae Graecae* (TLG).

Outcome: The first large-scale, open-access OCR and silver linguistic resource for Ancient Greek with **6M words** (< 1.05% of CER).

CHALLENGES

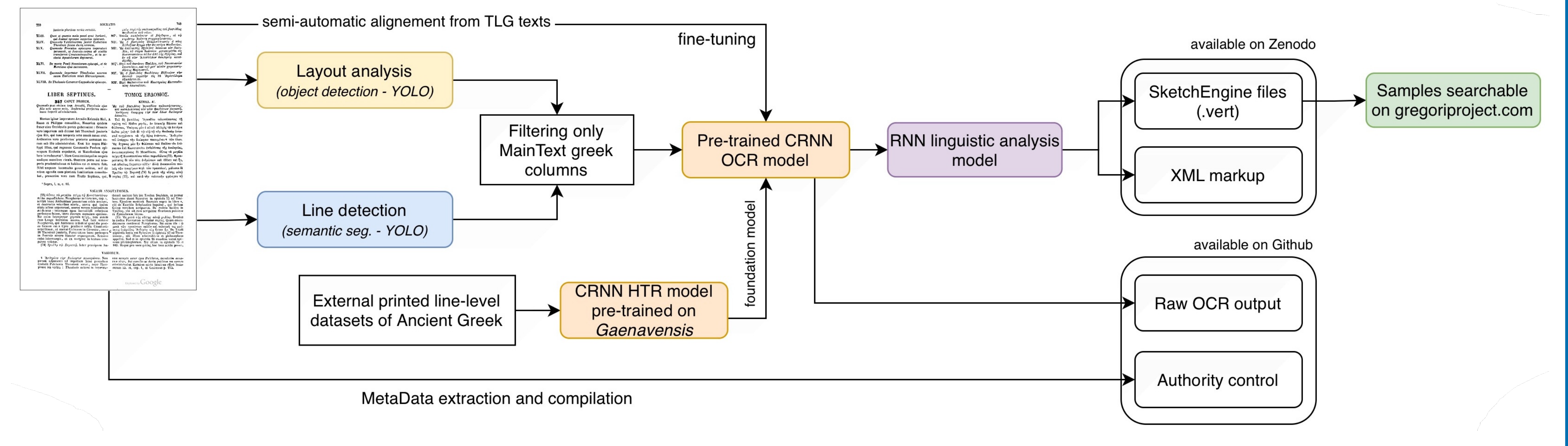


TEXT EXTRACTION AND TEXT ANALYSIS PIPELINE

To overcome the complex bilingual layouts (with crossing-lines) and degraded 19th-century typography, the pipeline relies on a **YOLO12s-detect model** for semantic zone detection and a **YOLO12s-seg model** for line segmentation. The text is then extracted using a **CRNN-based text recognition** (Vidal-Gorène *et al.*, 2021) pretrained with raw alignments from the TLG,

then fine-tuned for the *Patrologia Graeca*, achieving a 1.05% Character Error Rate (CER).

Following text extraction, the pipeline applies a hybrid strategy combining the PIE architecture (Manjavacas *et al.*, 2019), already pretrained with GREgORI, and a **dictionary-based post-correction** (Kindt *et al.*, 2022).



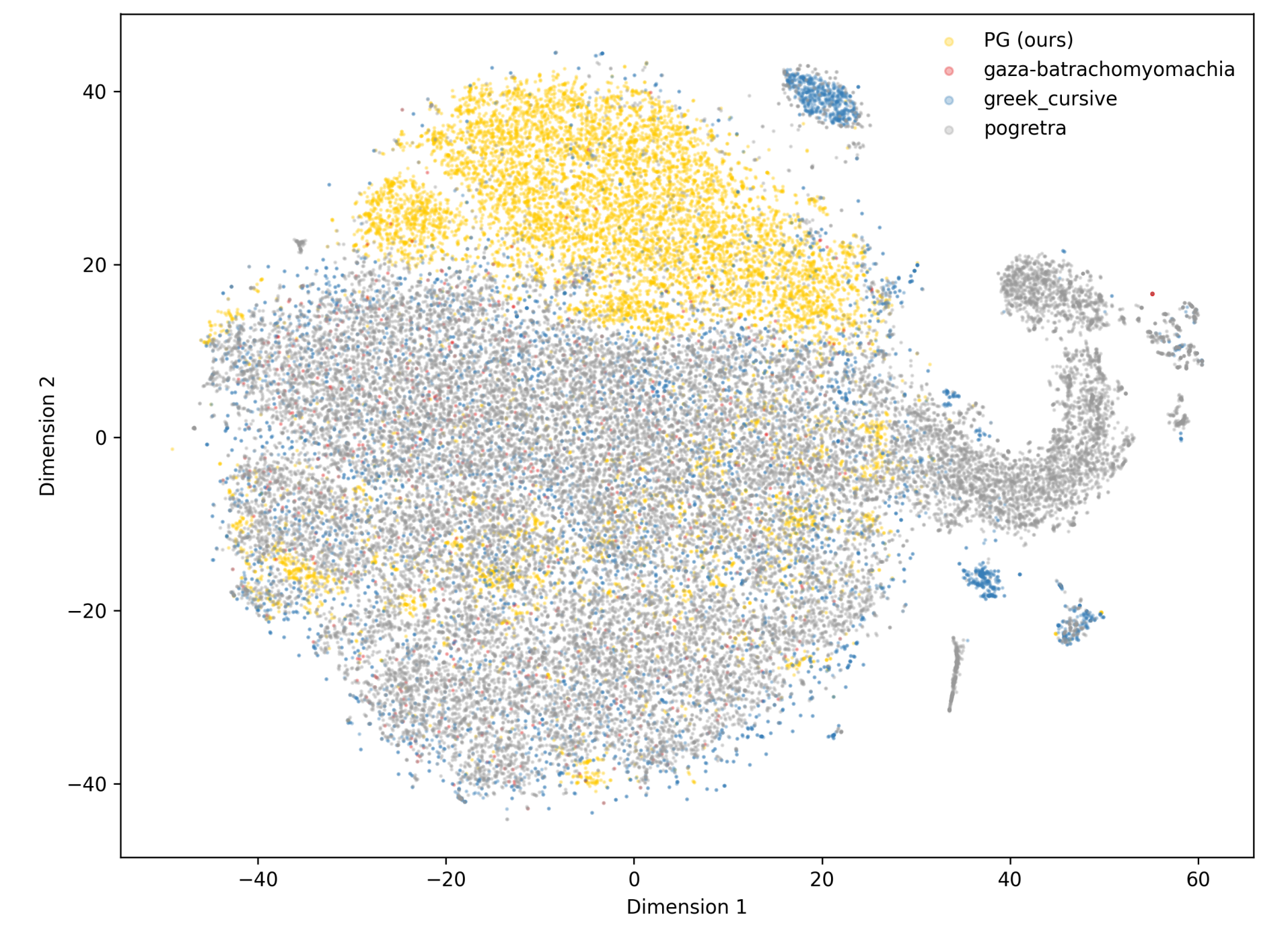
RESULTS AND CORPUS COVERAGE

Our fine-tuned OCR yields **1.05% CER** and **4.69% WER**, thanks to the **pretraining on real data** ($\approx 11k$ lines). Layout and line detection achieve high precision ($mAP50 \geq 0.97$). Notably, 80% of remaining errors are minor polytonic diacritic confusions rather than base character errors (no confusion observed at the lemmatization step). The resulting corpus comprises **6M lemmatized tokens** with full morphosyntactic annotation. t-SNE highlights the PG cluster (yellow) as a distinct semantic space driven by its specific vocabulary, a key for training new diachronically robust Ancient Greek language models.

Text recognition			
Model	CER	WER	
Tesseract	11.57%	39.65%	
Transkribus	6.14%	14.82%	
Ours (PG)	1.05%	4.69%	

Layout detection			
Class	P	R	mAP50
Greek column	0.963	0.994	0.973
Latin column	0.969	1.000	0.982
Title	0.462	0.950	0.462
Marginalia	0.422	0.891	0.427

Line detection			
Detection	0.983	0.994	0.973
Reading order	0.980	—	—



t-SNE: PG forms distinct semantic cluster, reflecting GT rich vocabulary compared to existing datasets.

Fragment of text from the Patrologia Graeca Corpus, showing the original Greek text and its transcription with diacritics.

Fragment of text from the Patrologia Graeca Corpus, showing the original Greek text and its transcription with diacritics.

RELEASE

Searchable corpus:

<https://www.gregoriproject.com>

Full raw data (OCR and tagged texts, .txt/.vert):

<https://zenodo.org/records/15780625>

OCR ground truth:

<https://github.com/calfa-co/Patrologia-Graeca>

Screenshot of the Patrologia Graeca Corpus search interface, showing search results for the word 'ΧΡΗΣΤΙΑΝΟΣ'.

Lemma search in the corpus, displaying concordances, analyzes, context, metadata and linked to lexical resources

- Kindt *et al.*, *Analyse automatique du grec ancien par réseau de neurones [...]*, BABELAO, Louvain (2022)
- Vidal-Gorène *et al.*, *A modular and automated annotation platform [...]*, ICDAR, Geneva (2021)
- Manjavacas *et al.*, *Improving Lemmatization of Non-Standard Languages [...]*, ACL, Minneapolis (2019)