

*Ὅσοι ἄνθρωποι, τοσαῦται γνῶμαι*

Harmonizing Guidelines for Handwritten Text Recognition of  
Ancient Greek

DH2025, Lisbon, July 14-18<sup>th</sup>, 2025

<https://dhtr25.anthologiagraeca.org/>



**DH2025**  
ACCESSIBILITY  
& CITIZENSHIP



Université   
de Montréal

  
ENS DE LYON Institut de recherche  
et d'histoire des textes

Mathilde Verstraete, Maxime Guénette,  
Matenia Vlachou, Marianne Reboul  
& Marcello Vitali-Rosati

# Motivation

- Growing number of **datasets** for Ancient Greek
- Expanding literature on HTR and **annotation** approaches
- Increasingly robust and generalizable **models**
- Clear need for **interoperable** practices and datasets  $\neq$  silos

# Textual Variability and Fluidity

Diacritics  
(accents &  
breathings)

Spelling

Abbreviations,  
*nomina sacra*,  
ligatures

Script Types  
(bookhand,  
cursive)

Allographs

Visual polysemy

Capitalisation  
practices

Spacing, Punctuation  
& other Auxiliary  
Signs

**BROAD**

Temporal,  
Geographical,  
Dialectal,

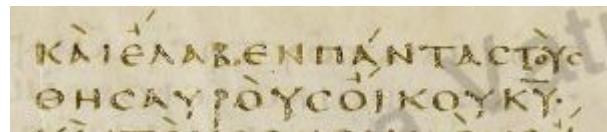
**DISTRIBUTION**

# Textual Variability and Fluidity

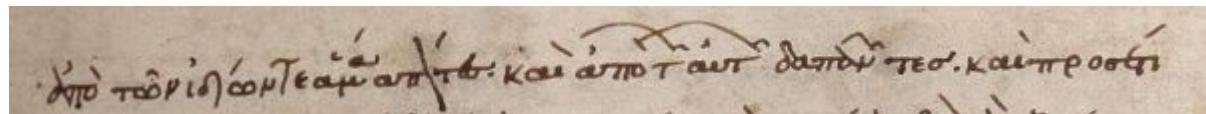
## Script Types and Capitalisation practices

From Majuscule  
to Minuscule

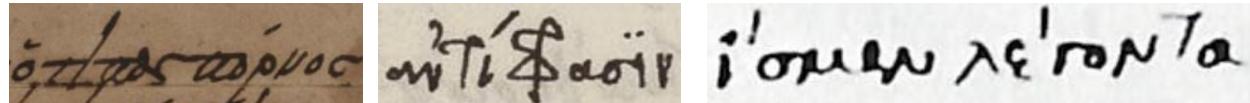
Various  
intermediate  
scripts



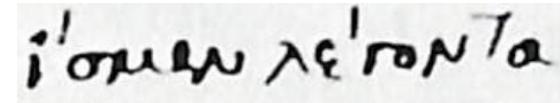
Vat. Gr. 1209 p. 426 (IVe s.)



Vat. Gr. 126, f. 33<sup>r</sup> (XIe s.)



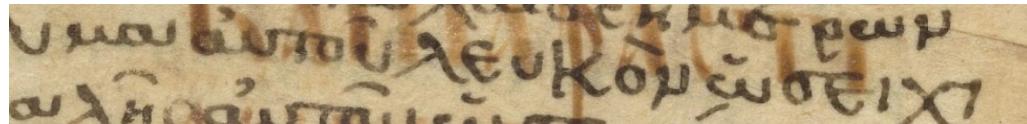
Vat. Gr. 2200, f. 53<sup>r</sup> (VIIe s.)  
Barb. gr. 221, f. 56v  
(XV-XVIe s.)



Ut. Arch. 258, f. 5<sup>r</sup> (XVe s.)

# Textual Variability and Fluidity

## Spacing, Punctuation & other Auxiliary Signs

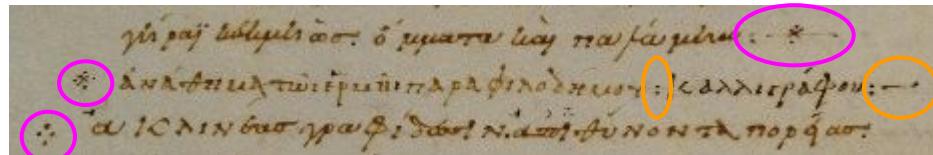
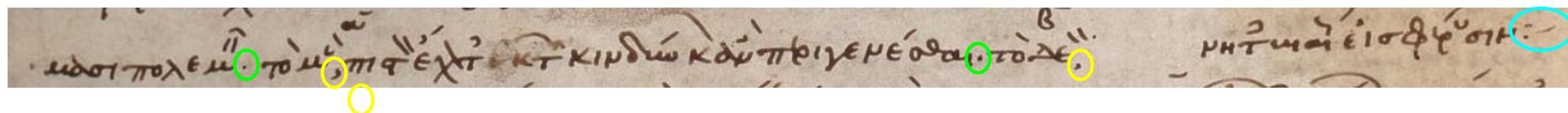


Bodl. Barocci 197\*, 2<sup>v</sup> (X-XIIIe s.)

ινουαύτοῦλευκομῶσειχι      (*scriptio continua*)

ἴνα οὗ αὐτοῦ λευκομῶσει χί      (separated)

Vat. gr. 126, f. 33<sup>r</sup> (XIe s.)



Pal. gr. 23, p. 153 (Xe s.)

# Textual Variability and Fluidity

## Spelling

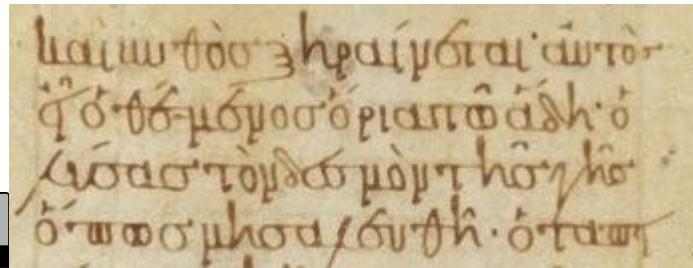
- [i] : ι η υ ει οι,
  - οῖκος or ἴκος,
  - χριστὸς or χρηστὸς,
- [o] : ο ου ω,
  - οῖκος or οῖκως,
  - λόγος or λώγος,
- [e] : ε or αι,
  - ἔτερος or ἔταιρος,
  - κενός or καινός,
- Other morphological variants
  - crasis: ἔάω or ἔῶ,
  - apocope: λέγει or λέγ'

# Textual Variability and Fluidity

## Diacritics

- τόνοι (combining or not)
  - rough breathing: ḥ (U+0314)
  - smooth breathing: ḥ (U+0313)
  - acute accent: ó (U+1FFD or 00B4)
  - grave accent: ò (U+1FEF or 0060)
  - **perispomeni:** õ (U+1FC0 or 0020)
- χρόνοι (combining or not)
  - breve: ɔ (U+0306)
  - **macron:** ᷑ (U+0304)

- διαστολή: , (U+002C) ?
- ύφεν: ὥ (U+032E) ?
- ἀπόστροφος: ’ (U+02BC) ?
- iota subscript: ὶ (U+0345)
- διαίρεσις: ὅ (U+00A8)



Par. gr. 1470, f. 135<sup>v</sup> (890).

Vat. Urb. gr. 33, f. 31r (s. XV).

A photograph of another page from an ancient Greek manuscript. This page appears to be a continuation of the previous one, featuring more of the same dark ink on aged paper. The script and diacritical marks are consistent with the first image, showing a variety of forms used throughout the text.

# Textual Variability and Fluidity

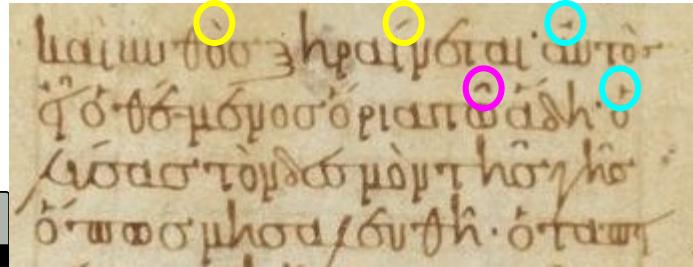
## Diacritics

- τόνοι (combining or not)

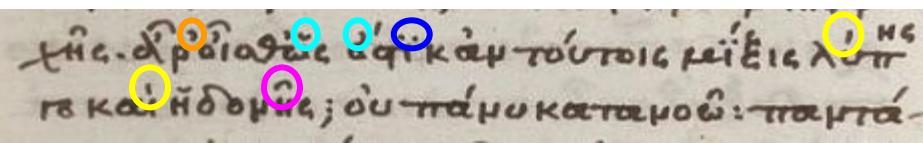
- rough breathing: ḥ (U+0314)
- smooth breathing: ḥ (U+0313)
- acute accent: ó (U+1FFD or 00B4)
- grave accent: ò (U+1FEF or 0060)
- perispomeni: õ (U+1FC0 or 0020)

- χρόνοι (combining or not)

- breve: ḡ (U+0306)
- macron: ḫ (U+0304)



Par. gr. 1470, f. 135<sup>v</sup> (890).



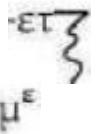
Vat. Urb. gr. 33, f. 31r (s. XV).

# Textual Variability and Fluidity

## Abbreviations, *nomina sacra*, ligatures

### Abbr. by suspension:

#### Desinences:

- -εται 
- -μενος 

#### Monograms:

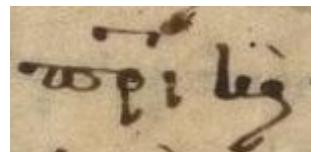
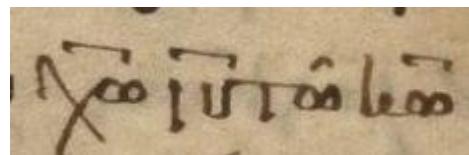
- μοναχός 

#### Numbers:

- X for χίλιοι
- H for ἑκατόν

### Abbr. by contraction:

- (Nomina Sacra)
- Sacred words,
- Proper names
- Numbers



Vat. gr. 2079, f. 109<sup>v</sup>

### Ligatures (e.g.):

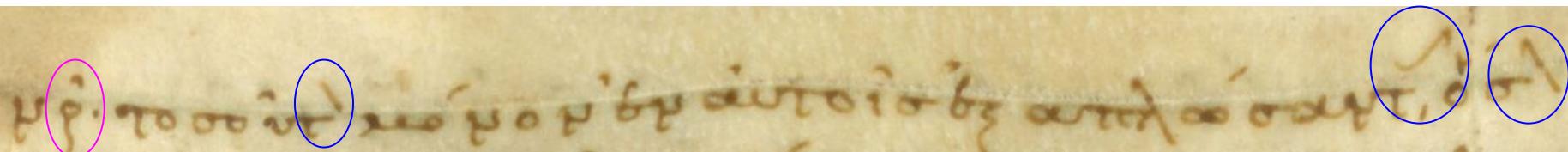


Agati (1984). Note paleografiche all' "Antologia Palatina". Boll. Class. 5, pp. 43-59 (fig.11)

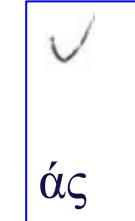
# Textual Variability and Fluidity

## Abbreviations, *nomina sacra*, ligatures

Vindob. Phil. G. 314, f. 110<sup>r</sup> (925)



νοῦν · τοσοῦτὸν μόνον ἐν αὐτοῖς ἐξαπλώσαντάς , δσον

v  · τοσοῦτ  μόνον ἐν αὐτοῖς ἐξαπλώσαντ ,  δσ 

οῦν  ὥν 

# Textual Variability and Fluidity

## Abbreviations, ligatures

desinenze	preposizioni congiunzioni avverbi	parole d'uso corrente
αι	ἀπό	ἀδελφός
αν	γάρ	χάρασται
ας	διὰ	ήμουν
αις	ἐπί	χέντρον
		χεφάλαιον
{ειν	καὶ	χοινόν
{ιν	ἢ	λέγειν
{ην	κατὰ	λόγος
εν	μὲν	μάρτυς
ες	μετά	πάντα
{ης	μήποτε	πόνος
{εις	θτι	σημείωσαι
	οῦν	σημειώσεων
μένος	πάλιν	σχόλιον
οις	παρὰ	τρόπος
ον	περὶ	φασὶ
ος	πρὸς	φησὶ
ου	πρότερον	χορός
ουν	πώποτε	χριστός
ους	πῶς	χρόνος
ων	ὑπέρ	ώρα
ως	χωρὶς	ώραιον

visible components → fully  
merged into ligature

στ ο τ τ σ τ π

αγ	γ	ον	θ θ θ
αθ	θ	οθ	θι θ
αν	ην	ου	γ γ γ
αρ	ερ	ους	γτ γτ
ασ	εσ	πε	πε
		ππ	ππ ππ
γγ	γ	πτ	πτ πτ πτ
		ρω	ρ ρ
εγ	γ γ γ	σα	σ α
ει	ει ει ει	σε	σ ε
ελ	ει	σθ	σ θ
εν	ει ει ει	σθ	σ θ
εξ	ει ει ει	σπ	σ σ σ
επ	ει ει ει	στ	σ τ τ τ
ετ	ει ει ει	σχ	σ χ
ευ	ει	τε	τ ε ε
		το	τ γ γ γ
εχ	ει	του	τ
εψ	ει	τρ	τ ξ ξ ξ
ην	ην ην ην	ττ	τ τ τ τ
ησ	ησ	τω	τ ω ω ω
ηθ	ηθ ηθ ηθ	υν	υ υ υ
ηε	ηε ηε ηε	υρ	υρ υρ
ητ	ητ ητ ητ	υσ	υς υς υς
ηλ	ηηηηηηηη	φρ	φ φ φ φ
ηο	ηηηηηηηη	ων	ω ω ω ω

E. Mioni (1973) *Introduzione alla Paleografia Greca*, Padova, p. 97 - 98

# Textual Variability and Fluidity

## Allographs and Morphological Evolution

sec. VI	sec. VII	sec. VIII	820/40
ααανανανα	αααανανα	αααανανα	αα
ββββδειη	ββββθειη	ββββθειη	β
γγγγγγγγ	γγγγγγγγ	γγγγγγγγ	γ
δδδδδδδδ	δδδδδδδδ	δδδδδδδδ	δ
εεεεεεεε	εεεεεεεε	εεεεεεεε	ε
ζζζζζζζζ	ζζζζζζζζ	ζζζζζζζζ	ζ
ηηηηηηηη	ηηηηηηηη	ηηηηηηηη	η
θθθθθθθθ	θθθθθθθθ	θθθθθθθθ	θ
ιιιιιιιι	ιιιιιιιι	ιιιιιιιι	ι
κκκκκκκκ	κκκκκκκκ	κκκκκκκκ	κ
λλλλλλλλ	λλλλλλλλ	λλλλλλλλ	λ
μμμμμμμμ	μμμμμμμμ	μμμμμμμμ	μ
νννννννν	νννννννν	νννννννν	ν
χχχχχχχχ	χχχχχχχχ	χχχχχχχχ	χ
οοοοοοοο	οοοοοοοο	οοοοοοοο	ο
ππππππππ	ππππππππ	ππππππππ	π
ρρρρρρρρ	ρρρρρρρρ	ρρρρρρρρ	ρ
σσσσσσσσ	σσσσσσσσ	σσσσσσσσ	σ
ττττττττ	ττττττττ	ττττττττ	τ
υυυυυυυυ	υυυυυυυυ	υυυυυυυυ	υ
φφφφφφφφ	φφφφφφφφ	φφφφφφφφ	φ
χχχχχχχχ	χχχχχχχχ	χχχχχχχχ	χ
ψψψψψψψψ	ψψψψψψψψ	ψψψψψψψψ	ψ
ωωωωωωωω	ωωωωωωωω	ωωωωωωωω	ω

σ

ε η θ ι

π

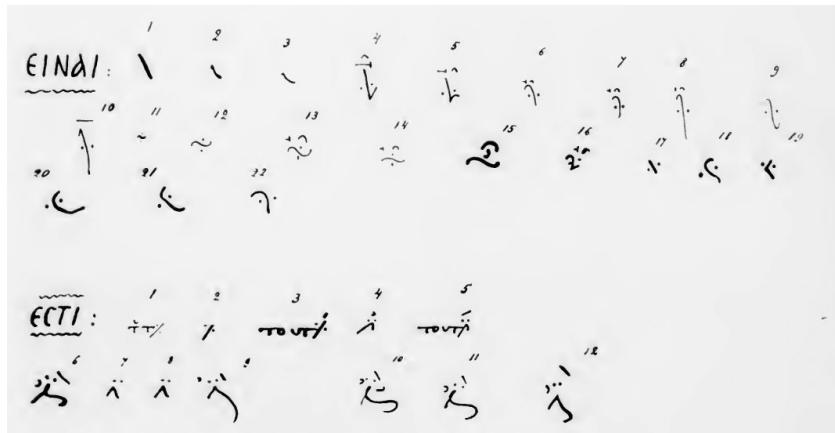
ρ π

β

υ β

# Textual Variability and Fluidity

## Visual polysemanism



// = -εισ, -ησ, -ισ

«Κη» = κη

# Transcribing (in theory)

A non-universal notion:

- There is no single "correct" transcription
- Every transcription reflects methodological choices
- Transcription is always project-driven

Transcription = Interpretation:

- Selection of what is deemed relevant
- Normalization according to conventions
- Contextualization within a framework of use

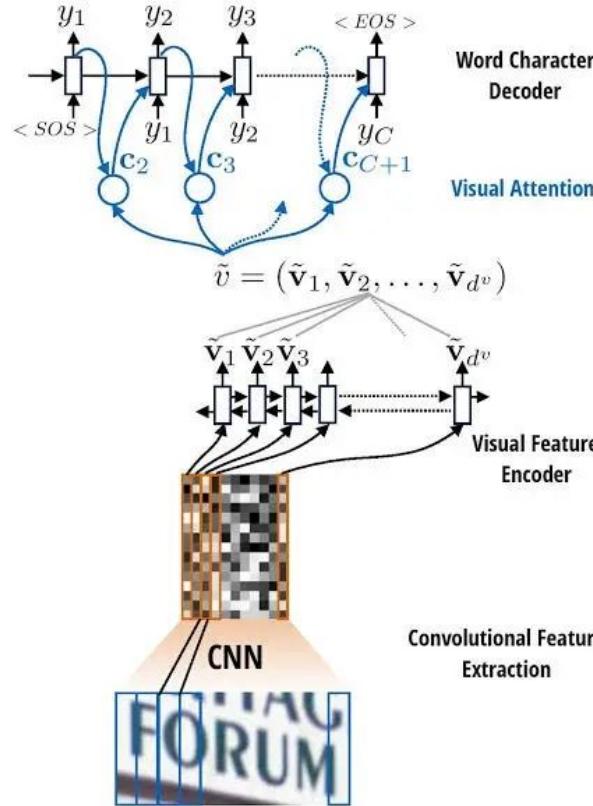
# Transcribing (in practice)

Limitations of Unicode... :

- **Incomplete coverage** of:
  - some ligatures, abbreviations, diacritics combinations;
    - Should we always decompose them? (e.g.: γ or και?)
    - What do we do for those that don't exist in Unicode? Invent Markup?
  - some auxiliary signs;
- **Various encoding options** for a same form:
  - e.g.: á =
    - U+1F04 (precomposed)
    - or U+03B1 + U+0313 + U+0301 (decomposed)
  - oxia vs tonos, e.g.: á = U-1F71 (Greek Small Letter Alpha with Oxia) or U-03AC (... with Tonos) ?

... Result as a **consistency nightmare** for HTR !

# Classic HTR (RNN/LSTM): Visual Mapping



What it does...

- Visual “recognition”: converts visual patterns (**glyphs**) into character sequences → close proximity between visual and semantic
- “Remembers” indirectly immediate context (preceding/following tokens)

What it doesn’t do...

- Semantic interpretation
- Automatic abbreviation expansion (should this be part of HTR?)
- Spelling correction or normalization (spacing, punctuation, capitalization)

# Limitations for Historical Material

## What HTR can do...

- Large-scale textual analysis
- Content discovery
- First-draft transcriptions

## What HTR cannot do...

- Perfect transcription
  - Precise paleographical analysis
- Need for post-correction

## Representativity → Allographs

- Rare variants underrepresented
- HTR struggles with unseen forms

# New age HTR (ViT/LMM): Beyond recognition

## Potential capabilities:

- Leverage external knowledge → contextualize and normalize
- Resolve abbreviations automatically
- Work with fewer examples (few-shot learning)

Base unit: ~~Glyph/Character~~ → Word-level or subword tokens

Limitations & challenges: Hallucinations (how can they be controlled?)

## Key questions:

- What should training data look like?
- Which tasks belong in HTR vs. post-processing?
- What task do we entrust the machine with?

# Transcription choices

Three\* main transcription approaches:

1. Graphetic
2. “Editorial”/”Regularized”
3. Graphemic
4. \* Mixed

→ The choice depends on the desired granularity of information

1.

ælo quando descendio ala tier  
cielo quando desçendio ala tier

2.



Passio sanctorum apostolorum Petri et Pauli.

1 Cum uenisset Paulus Romanum, conuenerunt ad eum omnes Iudei dicentes: Nostram fidem, in qua natus es, ipsam defende. non est enim iustum, ut cum sis Hebraeus ex Hebreis ueniens, gentium te magistrum iudices, et incircumcisorum defensor factus tu cum sis circumcisus, fidem circumcisionis evacuas. cum ergo Petrum uideris, suscipe contra eius doctrinam, quia

3.

ælo quando descendio ala tier  
cielo quando descendio a la tier

# Existing Datasets

Dataset	Century	Script	Size (char.)	Approach	Abbr.	Diacritics	HTR Model
<a href="#">Eparchos</a>	16th CE	Late byz. minuscule	116 894	<u>Graphemic</u>	Resolved	✓	✗
<a href="#">Méléagre</a>	10th CE	Byzantine minuscule	114 273	<u>Graphemic</u>	Unresolved	✓	✓
<a href="#">Zenon Papyri</a>	3rd BCE	Cursive	5 850	<u>Graphemic</u>	Unresolved	✗	✗
<a href="#">Stavronikita</a>	14-16th CE	Mixed minuscule	103 080	<u>Graphemic</u>	Resolved	✓	✗
<a href="#">HPGTR</a>	8-17th CE	Mixed minuscule	64 952	<u>Graphemic</u>	Resolved	✓	✗
<a href="#">Chrysostomicus I</a>	10-14th CE	Mixed minuscule	56 791	<u>Graphemic</u>	Resolved	✓	✓

# This workshop

## Today's Goals

- Build an interdisciplinary community around HTR for Ancient Greek, bringing together philologists, paleographers and HTR users (*aka, you*);
- **Identify the challenges** and needs of the community;
- Encourage participatory, use case-driven and **collaborative decision-making**;
- Develop **shared guidelines** for the transcription and encoding of Ancient Greek.

Website: <https://dhtr25.anthologiagraeca.org/>

Shared notes:

[https://docs.google.com/document/d/1S0xRZb1ImEBfKL6xPu7KmB9dSr15MSYZlFP\\_S8ZwR2A/edit?usp=sharing](https://docs.google.com/document/d/1S0xRZb1ImEBfKL6xPu7KmB9dSr15MSYZlFP_S8ZwR2A/edit?usp=sharing)

## Today's Overview

- [09h-09h20 : Introduction]
- 09:20-11:00 : Case Studies
  - [#1 : M. Verstraete & M. Guénette]
  - #2 : E. Perdiki
  - #3 : A. Jambé
  - #4 : C. Vidal-Gorène (online)
  - #5 : I. Marthot-Santaniello
- 11:00-11:15 : Coffee break
- 11:15-12:30 : Discussions around guidelines

