Ὅσοι ἄνθρωποι, τοσαῦται γνῶμαι
Harmonizing Guidelines for Handwritten Text Recognition of Ancient Greek
DH2025, Lisbon, July 14-18th, 2025

# HTR for Byzantine manuscripts

Recognising Ioannes Chrysostomus,
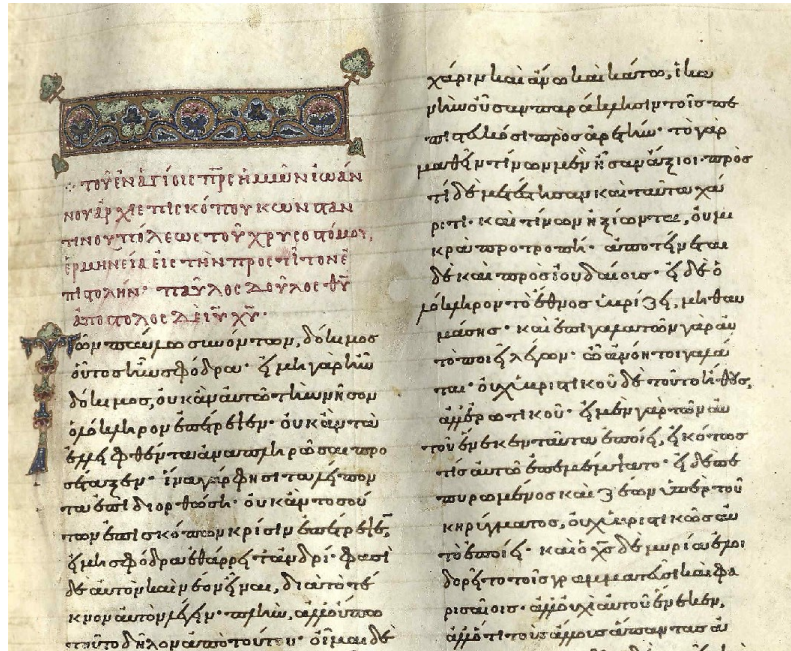Maximus Planudes, and Cyril of Alexandria
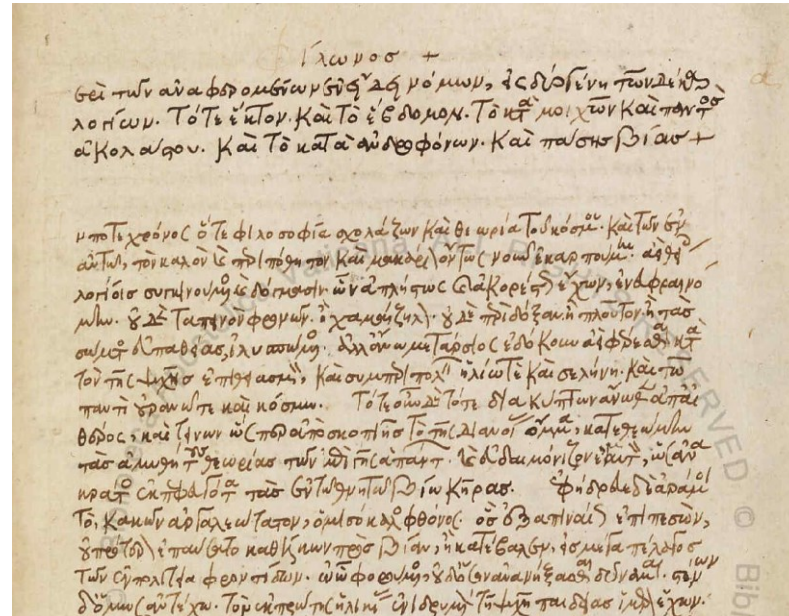
Elpida Perdiki
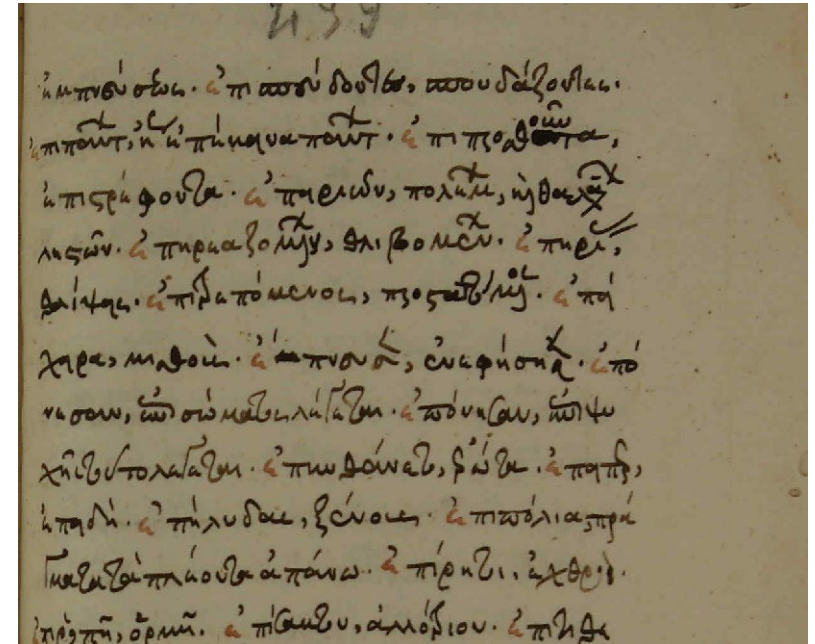Democritus University of Thrace

# What went wrong (and right)



John Chrysostom



Maximus Planudes



Cyril of Alexandria (?)

# One method, three case studies

Same system configurations for all models

Quality over quantity

No overlap between training & validation data

# General Workflow

- Digitised Images
- Diplomatic Transcription
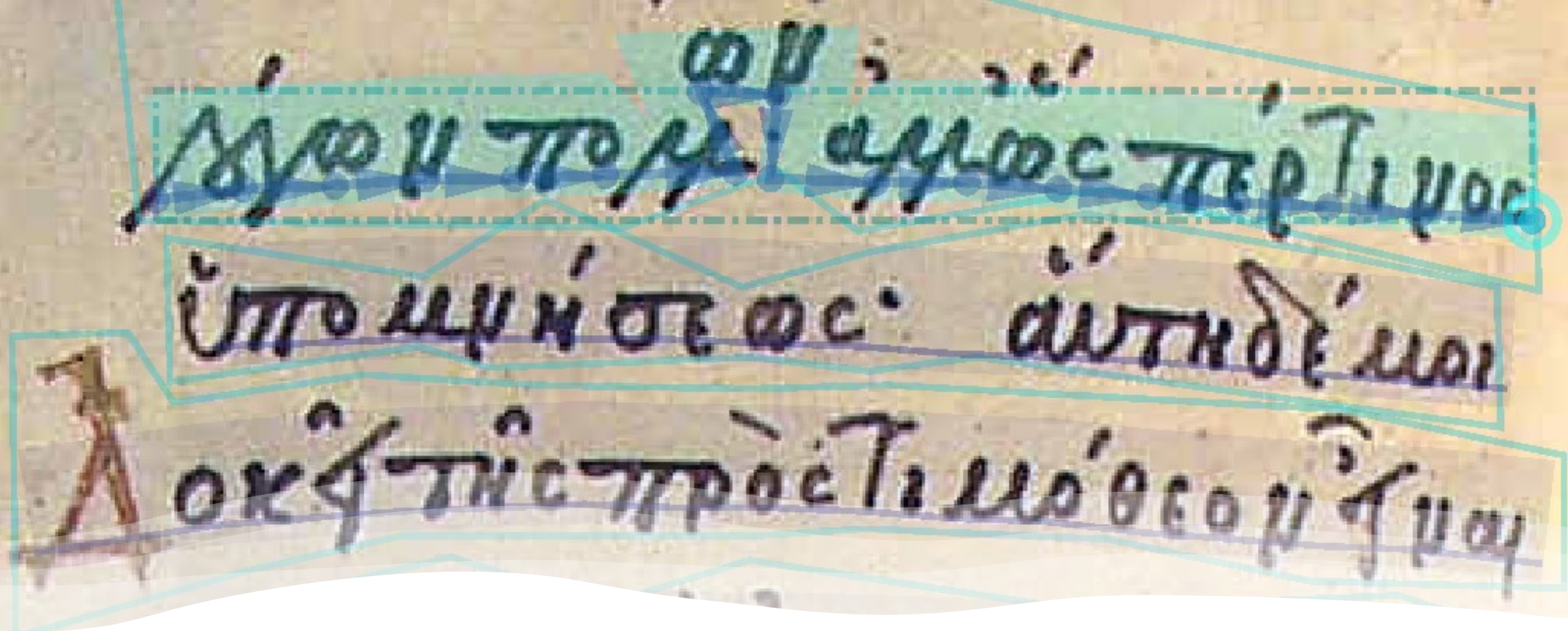- Model Training (Transkribus)
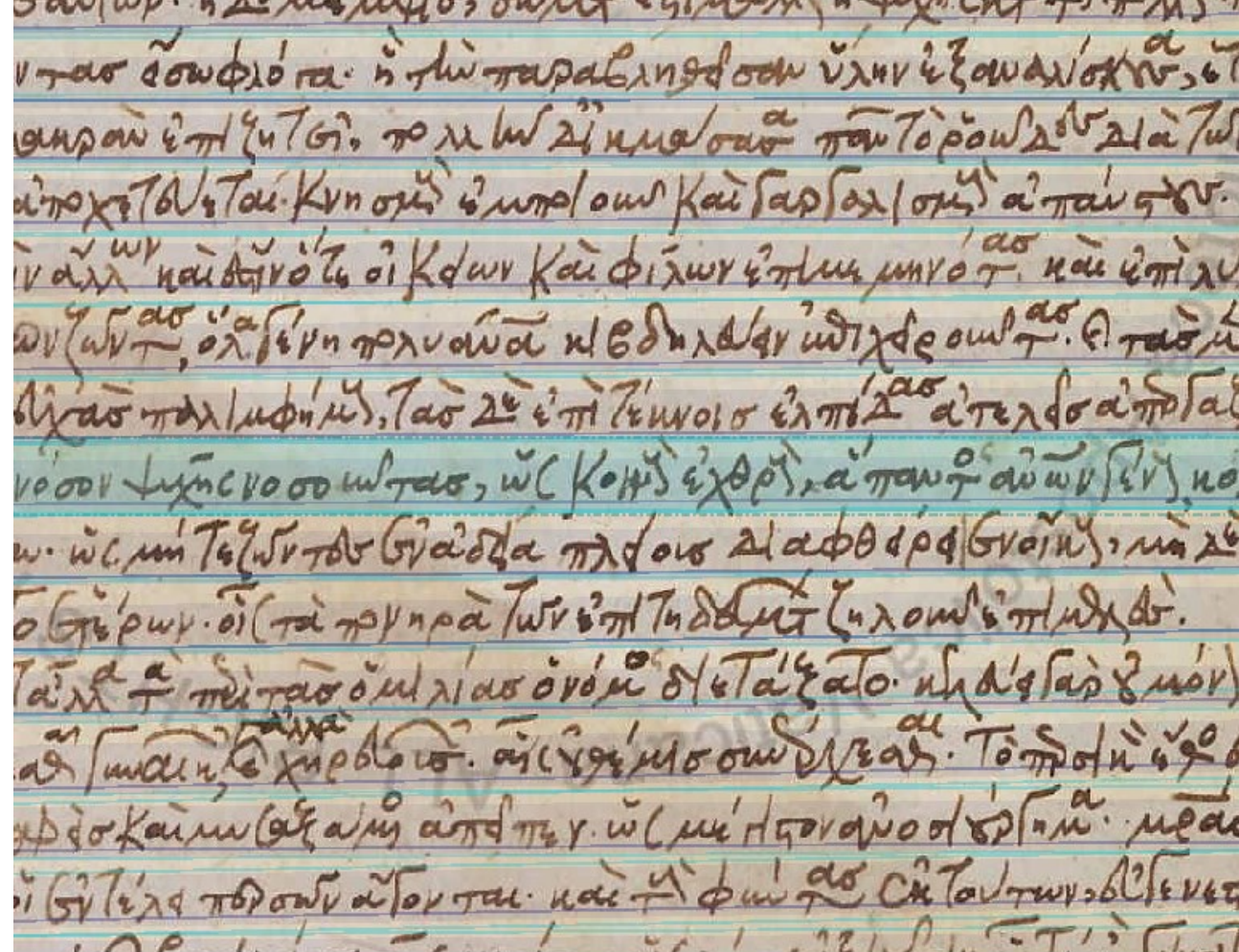- CER Evaluation
- Fine-tuning

## Case Study I: Ioannes Chrysostomus – The Great Success

- 11 manuscripts (10th–14th c.), mixed image quality
- Total training input: **29,228 words**
- Diplomatic, consistent transcription (2 annotators)
- Layout fine-tuned (superscripts, hanging lines)
- Final model (*Chrysostomicus I*): **3.9% CER**

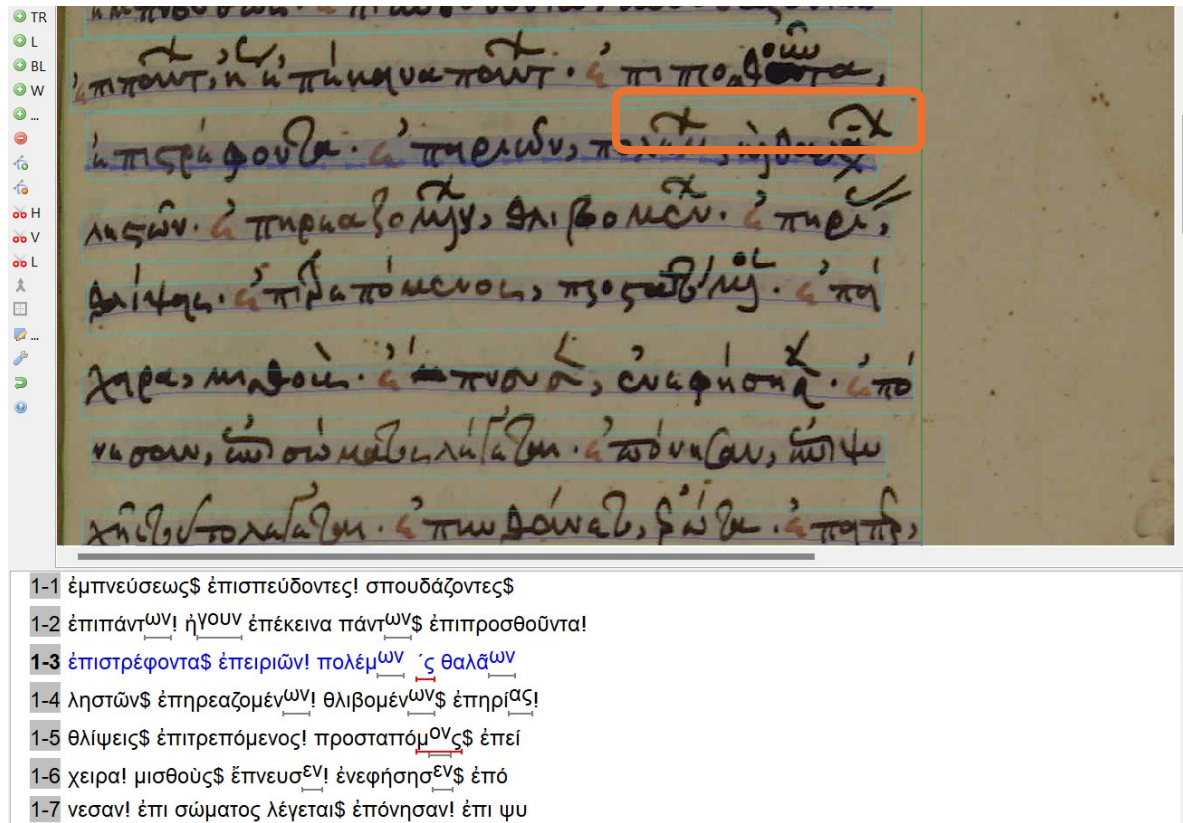# Maximus Planudes – The hidden potential of correcting your errors

- 2 autograph manuscripts, 13th century

- 4 transcribers, 8,000 words, diplomatic transcription

- Initial model (one manuscript): 16% CER

- After corrections: 8.5% CER

- Combined model (both manuscripts): 8.9% CER (transfer learning)

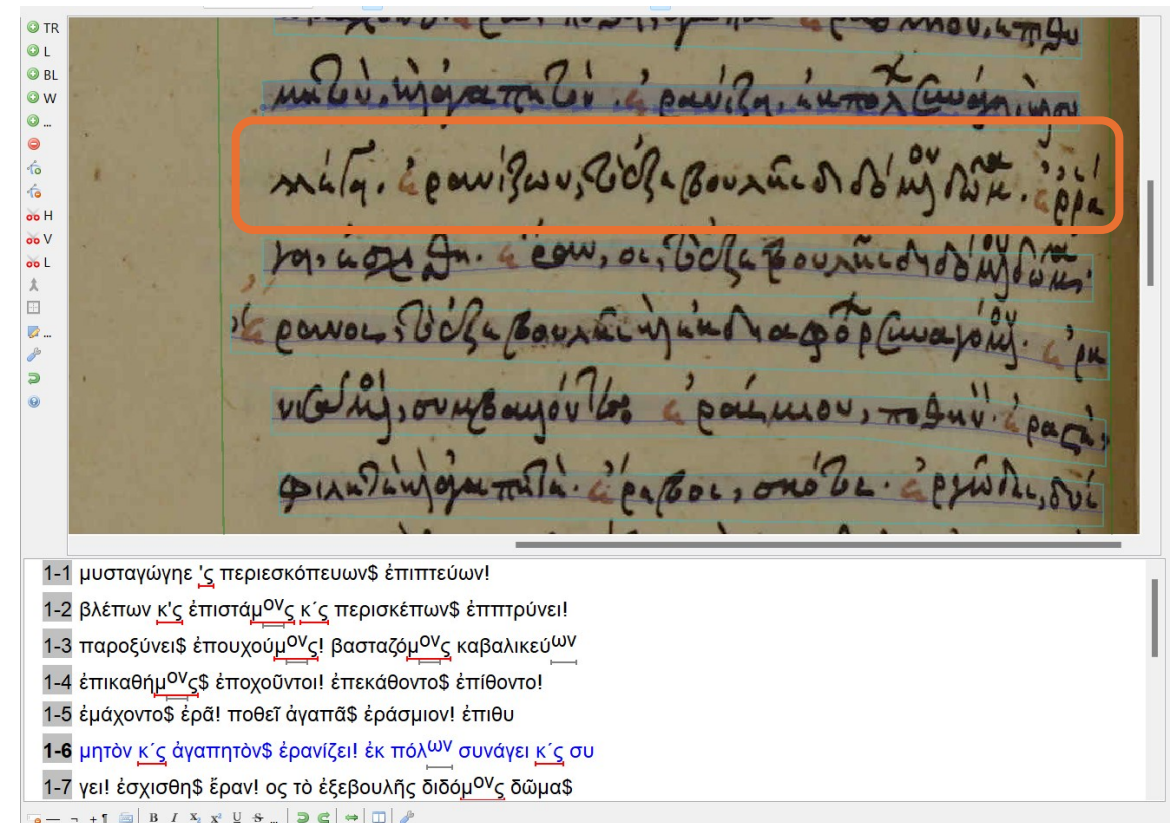- Main takeaway: transcription consistency > system tweaks



στιν ὅτε οἰκείων καὶ φίλων ἐπιμεμηνότ$^{ας}$, καὶ ἐπὶ λύμ$^η$ τῶν

γένη πολυάνᾶ κιβδηλεύειν ἐπιχειροῦντ$^{ας}$. χ᾿ τὰς μ$^{<᾿}$ ἐπιγα

ἡμ$^y$. τὰς δὲ ἐπὶ τέκνοις ἐλπίδ$^{ας}$ ἀτελεῖς ἀπεργαζομέν$^y$

οσοῦντας, ὡς κοιν$^y$᾿ ἐχθρυ᾿, ἄπαντ$^o$ ἀνῶν γέν$^y$ κολαστέ\

τες ἐν ἀδεία πλείους διαφθείρει ἐν οἴκ$^y$ μὴ δὲ διδάσκαλ$^{οι}$
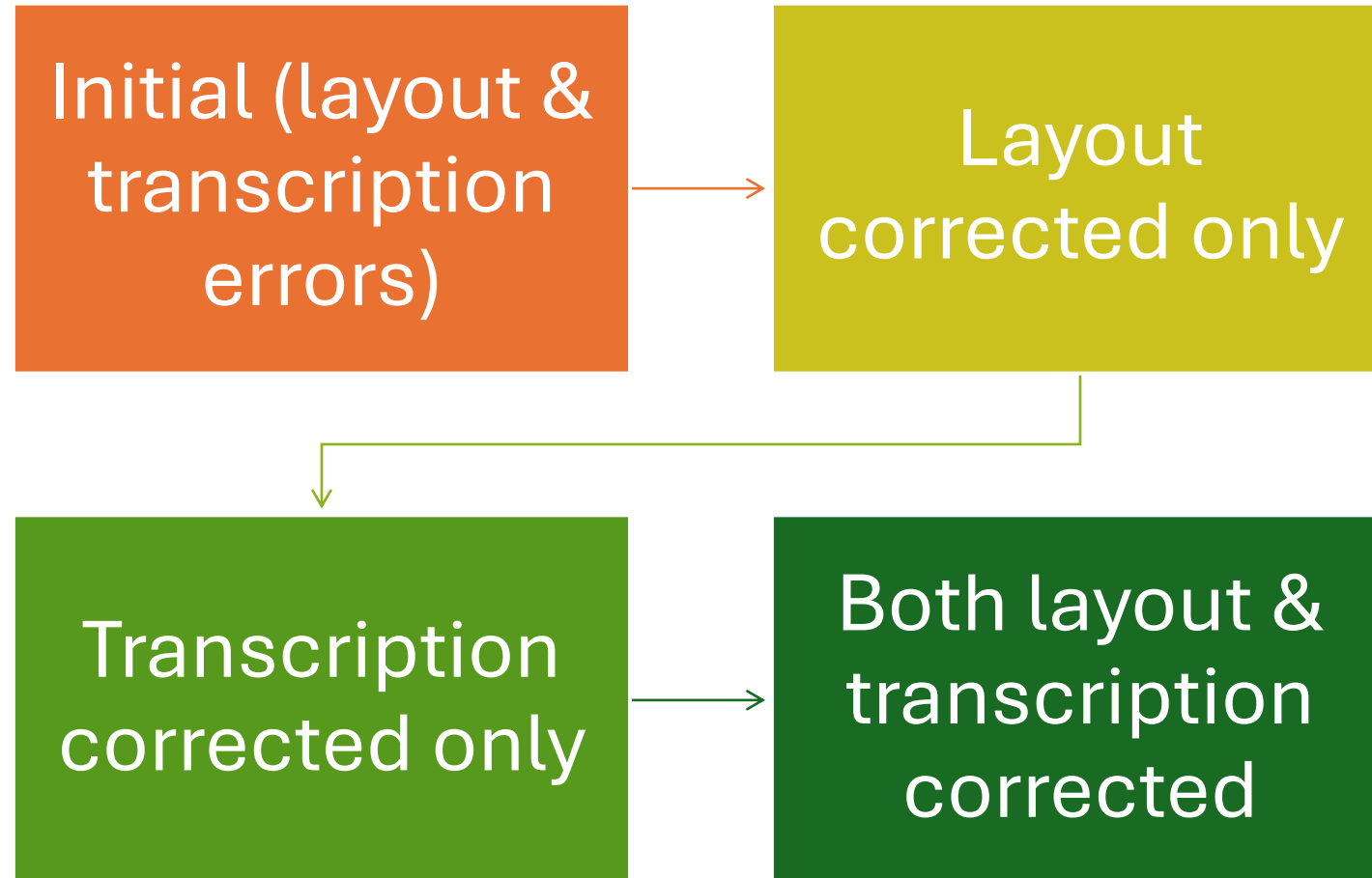
# Cyril of Alexandria – The almost disaster



Superscripted abbreviations out of basline

Missing line from layout recognition

# Targeting HTR errors: 4 different models

**Initial (layout & transcription errors)** → **Layout corrected only**

**Transcription corrected only** → **Both layout & transcription corrected**
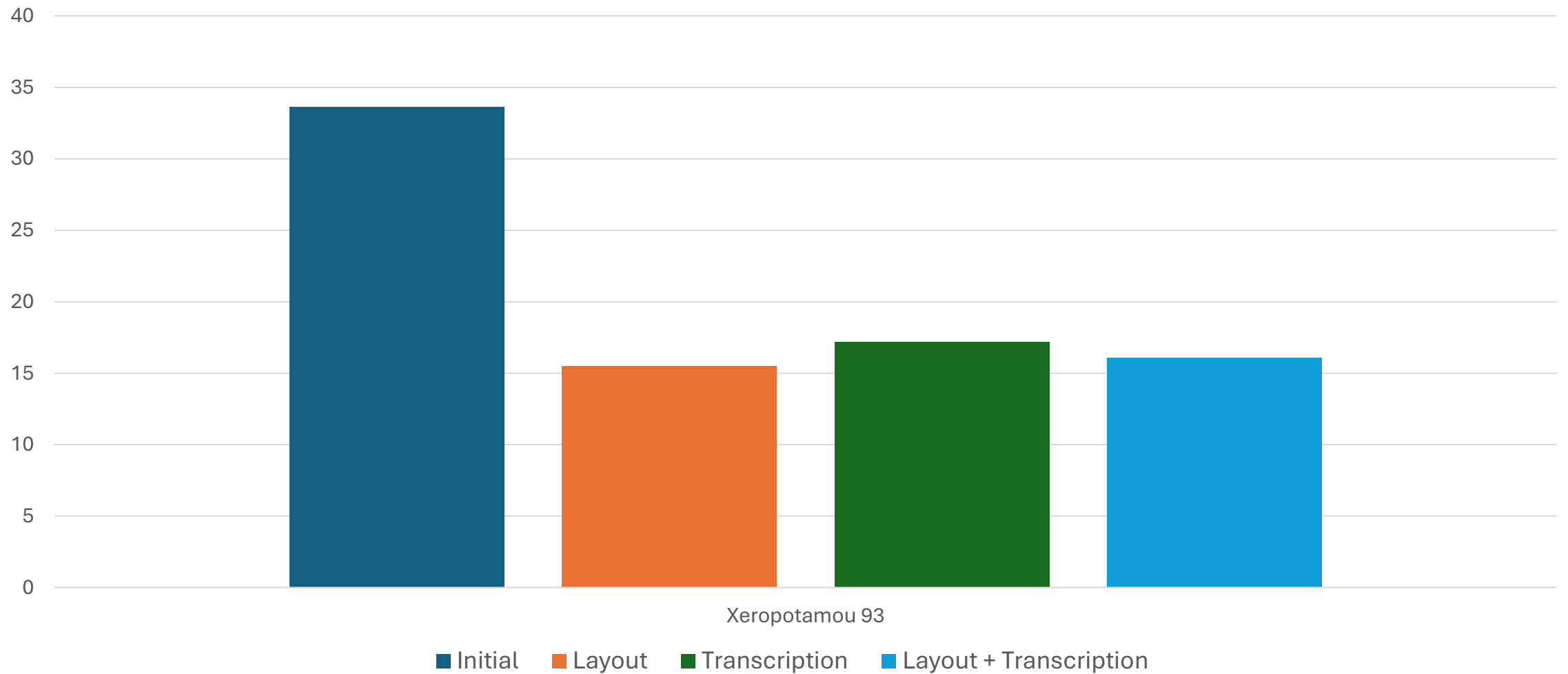
**Corrections applied only to:** Xeropotamou 093
**Training set:** 8,000 words (10% validation)
**System settings fixed:** max epochs=250, learning rate=0.0003, early stopping=20

# What it means

Minor corrections have large impact on performance

Layout correction yielded best CER drop

Transcription quality is equally critical

Correcting both offers most robust and stable model

# What we get from all of these

- HTR systems **depend heavily** on quality training data.
- Even clean digitisation **fails** if layout or transcription are inconsistent.
- **Quality & consistency** in annotation > quantity of data.
- Clear guidelines + supervision = **better performance**.
- The human factor is **as crucial** as technology.

# What needs to be done
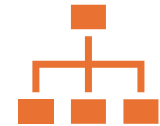
Systematic layout spot-checks before training

Peer-reviewed transcription tests & agreed conventions

Consistent abbreviation expansion (preferably in metadata)

Annotator training on guidelines & rationale

Ongoing communication & supervision

# All we have to decide is what to do with the time that is given us.

J.R.R. Tolkien, *The Fellowship of the Ring*