Prompting the Muse: Generating Prosodically-Correct Latin Speech with Large Language Models

Anonymous ACL submission

Abstract

This paper presents a workflow that compels an audio-enabled large language model to recite Latin poetry with metrically accurate stress. One hundred hexameters from the Aeneid and the opening elegiac epistula of Ovid's Heroides constitute the test bed, drawn from the Pedecerto XML corpus, where ictic syllables are marked. A preprocessing pipeline syllabifies each line, converts alien graphemes into approximate English-Italian counterparts, merges obligatory elisions, adds commas on caesurae, upper-cases every ictic syllable, and places a grave accent on its vowel. Verses are then supplied, one at a time, to an LLM-based Text-to-Speech model under a compact system prompt that instructs slow, articulated delivery. From ten stochastic realisations per verse, a team of Latin experts retained the best; at least one fully correct file was found for 91% of the 200 lines. Upper-casing plus accent marking proved the strongest cue, while hyphenating syllables offered no benefit. Remaining errors cluster around cognates where the model inherits a Romance or English stress template. The corpus of validated audio and all scripts are openly released on Zenodo, opening avenues for pedagogy, accessibility, and prosodic research.

1 Introduction

002

004

005

011

012

017

019

022

040

042

043

Latin prosody, at its core, is the systematic study of Latin poetry, particularly its laws of meter. Unlike English poetry, which relies on the alternation of stressed and unstressed syllables to create rhythm, classical Latin meter operates on a quantitative rhythm, determined by the arrangement of long and short syllables. The very term "prosody" finds its origins in the Greek word *prosoidia*, which initially signified a song sung to music or the specific pronunciation of a syllable.

Whereas handbooks faithfully describe reconstructed prosodical pronunciations, convincing spoken renditions accessible to learners remain scarce. Neural text-to-speech has closed the quality gap for modern languages, yet Latin remains marginal: the models lack training data and frequently transplant English or Romance stress patterns. 044

045

046

047

051

060

061

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Recent work in prosody editing offers an alternative. FastSpeech-type architectures expose duration, pitch, and energy predictors that can be edited after inference (Ren et al., 2020; Lam et al., 2025). Large language models with direct audio decoders add the possibility of steering pronunciation through plain text prompts, avoiding re-training. Their potential for historical languages has scarcely been explored.

The present study therefore asks whether prompt engineering, reinforced by symbolic prosodic annotation, is enough to make a general-purpose LLM read Latin verse with metrically correct stress.

2 Theoretical Background

2.1 Latin prosody

Classical verse rests on the opposition of long and short syllables, organised into metrical feet and regulated by fixed caesural patterns (Fortson IV, 2011). Quantity derives from vowel length and from consonantal environment, yet several phenomena blur the rule set: *muta cum liquida* allows optional resolution, while pervasive elision removes entire syllables at morpheme borders. Quantitative rhythm therefore resists categorical annotation; even the primary grammarians disagree in boundary cases. Because no contemporary acoustic evidence survives, phonological reconstruction must triangulate between Roman orthography, comparative Romance data, metrical practice, and prescriptive grammars (Allen, 1989). In practice, full reconstruction of absolute vowel length remains tentative. Modern pedagogy often replaces quantity with stress-based recitation, although stress in Latin is governed by its own moraic calculus. Any synthetic-speech system must decide which of

087

100

101

102

105

106

107

109

110

111

112

113

114

115

116

117

118

119

121

122

123

124

125

126

127

129

130

131

these competing principles to privilege.

2.2 Digital Latin: corpora, annotation, and prosodic tooling

Over three decades, Latin has moved from an almost text-only digital presence to a language with a modest but growing NLP stack (Riemenschneider and Frank, 2023). Tokenisers, lemmatisers, and treebanks are available through resources such as CLTK (Johnson et al., 2021), Stanza (Qi et al., 2020), and the Universal Dependencies Latin collections (De Marneffe et al., 2021). Prosodic annotation, however, remains rarer. Pedecerto annotates circa 244,000 dactylic lines from Musisque Deoque, returning syllabification, quantity, foot structure, and caesurae (Colombi et al., 2011). Its XML export supplied the gold data used in the present study. Other scanners address particular metres: the CLTK modules for hexameter and hendecasyllable (Johnson et al., 2021), Anceps for trimeters (Fedchin et al., 2022), and Loquax for quantitative syllabification and IPA transliteration (Court, 2025).

2.3 Large Language Models and **Prompt-Based Prosody**

Large language models trained on audio-text pairs have begun to encode prosodic regularities that can be elicited by prompt design. VALL-E and ZM-Text-TTS exploit massive multilingual corpora; their output retains speaker identity and sentence melody yet shows limited control over metre (Lam et al., 2025). The innovation proposed here inverts the usual pipeline: instead of sampling latent style tokens, we preprocess the poetic text, marking ictic positions and supplying approximate phonology in an orthography already mastered by the model (chiefly English with occasional Italian spellings for /u/ and palatals). At synthesis time those stress markers override default duration predictors, favouring long phones in ictic slots and shortened ones elsewhere. This approach follows the philosophy of PRESENT—prosody is steered through the input representation, not through additional parameters-yet applies it to classical verse rather than conversational prose.

2.4 Pedagogical and inclusive perspectives

Audio renditions of Latin verse remain an expen-128 sive commodity, created by a handful of trained classicists. Automated generation promises open collections usable in language teaching, literary

analysis, and accessibility contexts. Recent surveys in Digital Humanities stress the need for sharable, standardised, and FAIR corpora of recitations (De Sisto et al., 2024). By leveraging TTS engines and releasing the aligned text-audio pairs, the project aims to partially answer that call. Moreover, directing attention to stress rather than absolute quantity lowers the entry barrier for learners whose first language lacks phonemic length, while retaining a recognisable metrical pulse, in accordance with teaching standards across the world.

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

161

162

163

164

165

166

167

169

170

171

172

173

174

175

176

177

178

179

180

Methodology 3

3.1 **Corpora and metrical annotation**

The experiments draw on two well-known Latin texts: the opening one hundred hexameter lines of Vergil's Aeneid and the first elegiac epistula of Ovid's Heroides. Together they furnish examples of the two metres most frequently met in both school curricula and introductory prosody courses. A dactylic hexameter line consists of six feet, each prototypically realised as a long-shortshort (dactyl, D) or long-long (spondee, S) sequence; the fifth foot is normally a dactyl and the sixth is a spondee whose final syllable is anceps. The elegiac couplet pairs such a hexameter with a dactylic pentameter, divided by a diaeresis after the third arsis; in practice the pentameter is felt as two hemiepes with obligatory caesura.

Machine-readable scansion came from the Pedecerto project (Colombi et al., 2011). Pedecerto encodes each word with a sy attribute that enumerates syllables and marks ictic positions with an upper-case A. A fragment of the XML illustrates the structure:

```
line name="1" meter="H" pattern="DDSS">
  <word sy="1A1b" wb="CF">Arma</word>
  <word sy="1c2A2b" wb="CF">uirumque</word>
</line>
```

During import the parser retained verse boundaries, foot patterns, ictus markers, word-boundary flags, and elision hints.

3.2 Text preparation pipeline

Each line was passed through an iterative preprocessing routine and immediately spoken by a synthesis model; Latinists then annotated pronunciation errors, after which the routine was adjusted. Syllabification relied on the Classical Language Toolkit, whose rule-based engine already covers

enclitics and diphthongs (Johnson et al., 2021). A 181 grave accent was placed over the vowel of every ic-182 tic syllable and the entire syllable was upper-cased. Words forming obligatory elision were merged (quoque et \rightarrow quoquet) in accordance with the Pedecerto wb attribute. Caesurae were rendered 186 by a comma, but only when the manuscript trans-187 mitted no other punctuation at that position; this decision proved particularly useful for pentameter lines, where the pause after the third arsis is nearly 190 fixed. Trials in which syllables were separated by hyphens (ar-ma vi-rum-que) showed no measur-192 able benefit and were dropped. 193

194

195

196

198

199

201

204

206

207

210

211

212

213

214

215

216

217

218

219

222

227

228

231

Orthographic substitution aimed at a rough classical pronunciation that modern English or Italian acoustic models could approach. Stops before front vowels were written k instead of c; qu became kw; ae and oe became ai and oi; ge and gi were expanded to ghe and ghi.

Because long contexts tended to blur prosodic control, each verse was spoken in isolation. A verse forms a minimal rhythmic unit whose internal pattern must remain coherent, whereas inter-verse junctures tolerate short pauses.

3.3 Speech synthesis experiments

Two families of systems were compared. Conventional sequence-to-sequence TTS engines (Tacotron 2, Kokoro, tts-1, tts-1-hd) could not ingest elaborate instructions; their output misstressed Latin loans that resemble high-frequency English forms and showed erratic vowel quantity. Large language models with integrated audio decoders performed better, presumably because the system prompt can impose prosodic policy. Several models in the GPT-40 and Gemini lines were tested; gpt-40-mini-tts (Hurst et al., 2024) delivered the most consistent rhythm and segmental clarity.

Prompt engineering proceeded from a verbose style sheet to a compact directive. Lengthy system prompts improved intonational contour but occasionally confused stress placement. The final prompt retained only three imperatives: speak slowly, articulate every syllable, obey the marked stresses. Repeating the fully processed verse inside the prompt, exactly as the model should pronounce it, brought an unexpected improvement, perhaps because the acoustic decoder aligns its plan with the visible text.

As LLMs incorporate stochastic sampling, pronunciation varies across runs. For each verse ten realisations were generated. When specialists reviewed the set, at least one rendition met the acceptance threshold in 91 percent of lines. Most remaining errors involved lexical interference from Romance or English cognates; for instance, the word passus from the Aeneid's fifth line emerged as pàssus rather than the required passùs. Re-spelling the stressed vowel (passùus) in the prompt usually resolved the problem, though this fix was applied sparingly, since excessive vowel doubling sometimes misled the model elsewhere in the line.

233

234

235

236

237

238

239

240

241

242

243

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

Sequences with dense stress, such as spondaic clusters, challenged the model, as did runs of elided vowels or complex consonant groups. These limitations are examined in Section 5.

3.4 Expert evaluation protocol

Three scholars of Latin phonology, none involved in system development, evaluated every candidate recording. Errors were marked on a span basis and classified as segmental, stress, elision, or pacing. Feedback was returned after each experimental cycle, leading to successive refinements of preprocessing and prompts until the acceptance rate stabilised.

3.5 Dissemination of audio material

The highest-ranked file for each verse was retained. Verses were concatenated with 800 ms silences, yielding two continuous recitations that mirror performance practice yet preserve per-line rhythmic autonomy. Waveform-level normalisation ensures homogeneous loudness. The corpus, its preprocessing scripts, and annotation spreadsheets will be deposited on Zenodo; a DOI will be included in the camera-ready version.

4 Results

4.1 Quantitative assessment

The evaluation covered 216 autonomous lines, of which 158 hexameters and 58 pentameters. Ten recordings were generated for every line, yielding two thousand candidate files. Table 1 reports acceptance rates after expert screening. The final system prompt is as follows:

This is a Latin poetical verse. Pronounce it rhythmically, slowly and with emphasis, articulating each syllable and correctly stressing them. Pronounce it like this: [pre-processed verse]

Metre	Lines	Lines with at least
		one correct realisa-
		tion
Hexameter	158	91.1%
Pentameter	58	91.4%
Total	216	91.2%

Table 1: Overview of the obtained Latin verse record-ings.

280Incorrect spans fell into four categories: segmental281substitutions, misplaced ictus, elision failure, and282pacing anomalies. Inter-annotator agreement on283the five-way label reached $\kappa = 0.79$ for hexameter284and $\kappa = 0.84$ for pentameter.

4.2 Effect of preprocessing variants

285

286

287

291

292

293

296

297

309

310

311

312

313

Ablation tests, run on a ten-line subset to contain annotation effort, show that three operations account for most of the gain over a plain graphemic baseline:

- Upper-casing and accenting the ictic syllable considerably reduced stress errors;
- Orthographic substitution of c/qu/ae/oe and palatal stops diminished segmental errors;
- Explicit commas on caesura lowered pacing mistakes, especially in pentameters.

Conversely, syllable hyphenation had negligible impact, while long system prompts improved intonation without improving segmental or stress accuracy. These findings corroborate earlier observations by Lam et al. (2025) that explicit duration–pitch instructions dominate hidden stylistic embeddings in LLM-based TTS.

4.3 Listening quality

Mean opinion scores were collected from fourteen external listeners familiar with Latin recitation but naïve to the study. They judged naturalness and metrical fidelity on a five-point scale. Best-of-ten selection reached 4.1 ± 0.6 for hexameter and 3.9 ± 0.7 for pentameter. Ratings drop by roughly one point when a random sample rather than the best file is played, reflecting the intrinsic variance of stochastic decoding.

5 Conclusions and outlook

The workflow demonstrates that a contemporary audio-enabled large language model, guided by minimal yet well-targeted textual cues, can read classical Latin verse with a promising degree of prosodic correctness. Stress salience carried by case-shift and diacritic proved a stronger cue than any attempt at modelling moraic weight directly, an outcome consistent with linguistic evidence on the rhythmical importance of stress in Latin poetry (Pawlowski and Eder, 2001). Segmental confusion arises chiefly from orthographic overlap with Italian and English; paradoxically, rare or morphologically opaque words are rendered more faithfully because no competing template exists in the model's training distribution. 316

317

318

319

320

321

322

323

324

325

326

327

328

329

331

332

333

334

337

338

340

341

342

343

344

346

347

348

349

350

351

352

353

354

355

356

358

5.1 Future work

Two lines of research appear promising. First, coupling the current prompt-based strategy with the controllable duration and energy interfaces available in FastSpeech-type decoders (Ren et al., 2021) may supply the missing quantitative layer. Second, training a lightweight alignment model on our validated recordings would allow deterministic selection rather than trial-and-error sampling. Beyond technology, the public release on Zenodo of both source XML and mastered audio will facilitate studies in metrics, second-language acquisition, and accessibility. The same pipeline applies, mutatis mutandis, to other Greco-Roman metres, to post-classical accentual hymns, and even to vernacular verse traditions where scholarly recordings are scarce.

Limitations

The system remains probabilistic. A user must be willing to request several readings and to curate the output manually. Dense spondaic passages, intricate elisions, and clusters such as ctn or gns still trigger mis-timed syllable nuclei. Quantity is approximated through pacing alone; true heavylight contrast, audible as durational ratio, is not yet guaranteed. Finally, the present study uses a single North-Atlantic vocal profile, whereas pedagogy would profit from multiple voices and speaking rates.

Acknowledgments

No further acknowledgements are declared.

4

References

361

362

363

367

373 374

375

381

387

391

394

395

396

400

401

402 403

404

406

407

408

409

410

411

412

413

414

- W Sidney Allen. 1989. Vox Latina: a guide to the pronunciation of classical Latin. Cambridge University Press.
- Emanuela Colombi, Luca Mondin, Luigi Tessarolo, Andrea Bacianini, Dylan Bovet, and Alessia Prontera.
 2011. Pedecerto. *Pedecerto. Metrica Latina Digitale*.
- Matthieu Court. 2025. Loquax: Nlp framework for phonology. https://github.com/mattlianje/loquax. GitHub repository.
- Marie-Catherine De Marneffe, Christopher D Manning, Joakim Nivre, and Daniel Zeman. 2021. Universal dependencies. *Computational linguistics*, 47(2):255– 308.
- Mirella De Sisto, Laura Hernández-Lorenzo, Javier De la Rosa, Salvador Ros, and Elena González-Blanco. 2024. Understanding poetry using natural language processing tools: a survey. *Digital Scholarship in the Humanities*, 39(2):500–521.
- Aleksandr Fedchin, Patrick J Burns, Pramit Chaudhuri, and Joseph P Dexter. 2022. Senecan trimeter and humanist tragedy. *American Journal of Philology*, 143(3):475–503.
- Benjamin W Fortson IV. 2011. Latin prosody and metrics. A companion to the Latin language, pages 92– 104.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Kyle P Johnson, Patrick J Burns, John Stewart, Todd Cook, Clément Besnier, and William JB Mattingly. 2021. The classical language toolkit: An nlp framework for pre-modern languages. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing: System demonstrations*, pages 20–29.
- Perry Lam, Huayun Zhang, Nancy F Chen, Berrak Sisman, and Dorien Herremans. 2025. Present: Zeroshot text-to-prosody control. *IEEE Signal Processing Letters*.
- Adam Pawlowski and Maciej Eder. 2001. Quantity or stress? sequential analysis of latin prosody. *Journal* of Quantitative Linguistics, 8(1):81–97.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D Manning. 2020. Stanza: A python natural language processing toolkit for many human languages. *arXiv preprint arXiv:2003.07082*.
- Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. 2020. Fastspeech
 2: Fast and high-quality end-to-end text to speech. arXiv preprint arXiv:2006.04558.

Frederick Riemenschneider and Anette Frank. 2023. Exploring large language models for classical philology. *arXiv preprint arXiv:2305.13698*.