# AI-Driven Approaches to Latin Language Modeling

## Data Driven Humanities Team*

## Background

The **Data Driven Humanities Research Group** is at the forefront of integrating artificial intelligence with classical philology. This poster presents the latest work of our group, showcasing our research in AI and machine learning for Latin language modeling. Our research focuses on developing a Latin language model by integrating lemmatization, sentence modeling, and prompt engineering, both through traditional querying and automated methods, to enhance linguistic processing, textual comprehension, and language learning.

## Latin Lemmatizer

We developed a Latin Lemmatizer that displays declension tables for any queried word(s). The basis of the lemmatizer comes from an online Latin dictionary website. For each word, webscraping is utilized to access every single table and its corresponding word forms. There are functions which detect the part of speech of a given word, the declensions of a table, and printing out all possible words that the website has access to. Right now, all parts of speech work (i.e., noun, verb, adjective, adverb, pronoun, preposition, conjunction, interjection), though there are some minor caveats: words with multiple possible origins (e.g., pulla). The algorithm right now throws an error when this happens, and it asks the user to specify the part of speech they are looking for. If the user re-queries with the desired word and part of speech, then it will choose the first result with that part of speech (see figure 1). Subsequently, we use the lemmatizer for sentence generation to pretrain our sentiment analysis model.

## Latin Sentiment Analysis Model

Our Latin Sentiment Analysis system is built upon the *"pnadel/LatinBERT"* model, which is a specialized BERT model pre-trained on Latin text. We modified this base model for sentiment analysis by:

1. Replacing the original masked language modeling head with a sequence classification head
2. Configuring the model for binary classification (positive/negative sentiment)
3. Implementing a custom tokenizer (*"pnadel/latin_tokenizer"*) specifically designed for Latin text
4. Adding padding token configuration to handle variable-length sentences

## PNADel/LatinBERT

### PNADel/LatinBERT (Hugging Face Model)

- **Developer:** Hugging Face user/org `PNADel`
- **Platform:** Fully trained and deployed via Hugging Face Transformers
- **Architecture: BERT-base** or similar, pretrained from scratch (not based on mBERT)
- **Corpus:** Likely a **larger, more heterogeneous corpus**, possibly including Classical, Medieval, and Ecclesiastical Latin
- **Purpose:** Broader applicability — usable for a wide variety of **downstream NLP tasks** in Latin via transfer learning



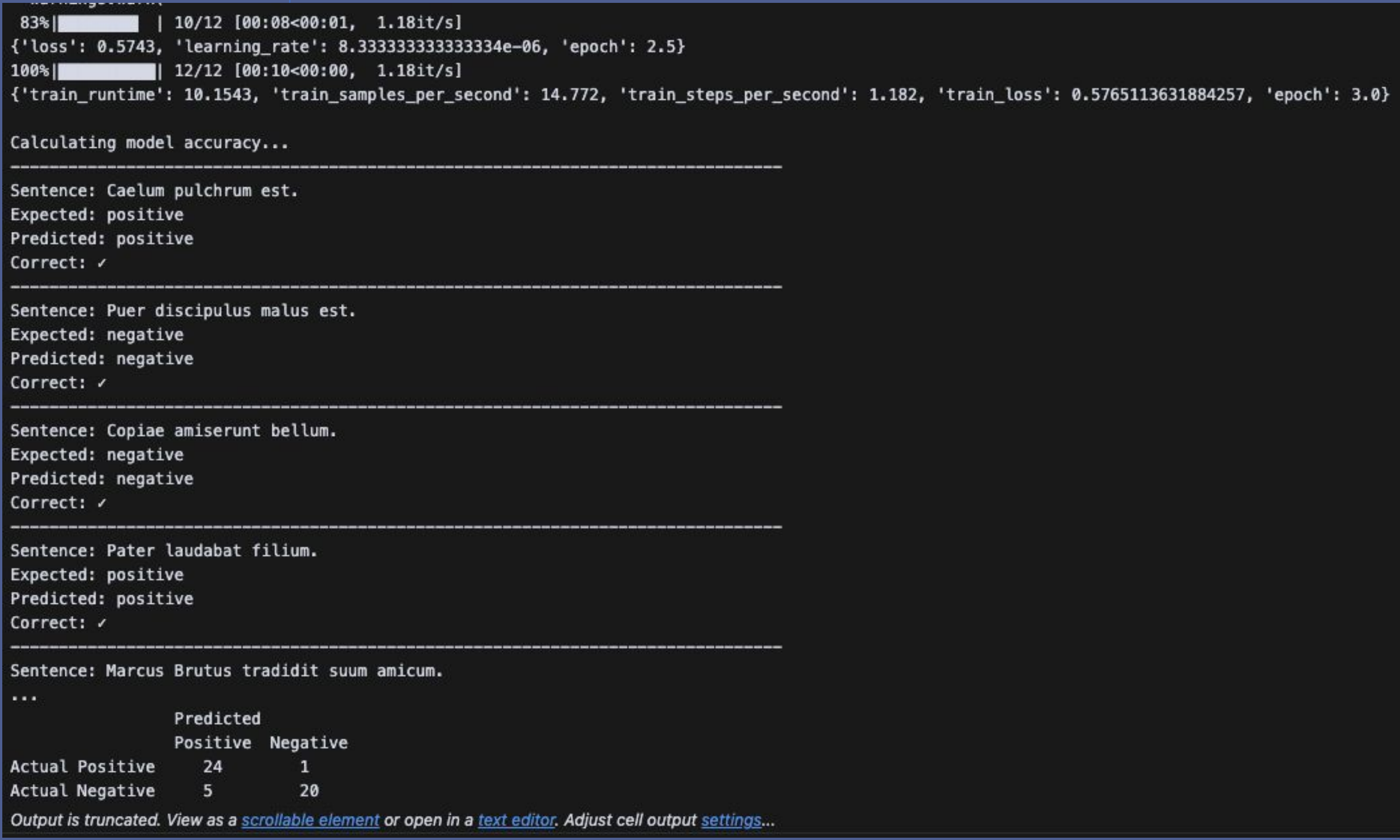**Figure 1:** Web response when querying "pulla"



**Figure 2.** Our modified model achieved an impressive 88% accuracy on the test set, demonstrating strong performance in classifying Latin sentences. The model successfully handles various Latin sentence structures and maintains consistent performance across different types of content.

## Query Engineering and Automation

Another focus of our project is the development and use of innovative methods in computational philology through **advanced query engineering** for Greek and Latin texts. Initially, our researchers manually crafted complex queries using ChatGPT-4 to analyze syntactic structures in ancient languages, allowing for precise linguistic and interpretive insights.

Building on this foundational work, the team has now programmed these processes into an automated system, training the model to recognize and process Greek and Latin syntax autonomously.
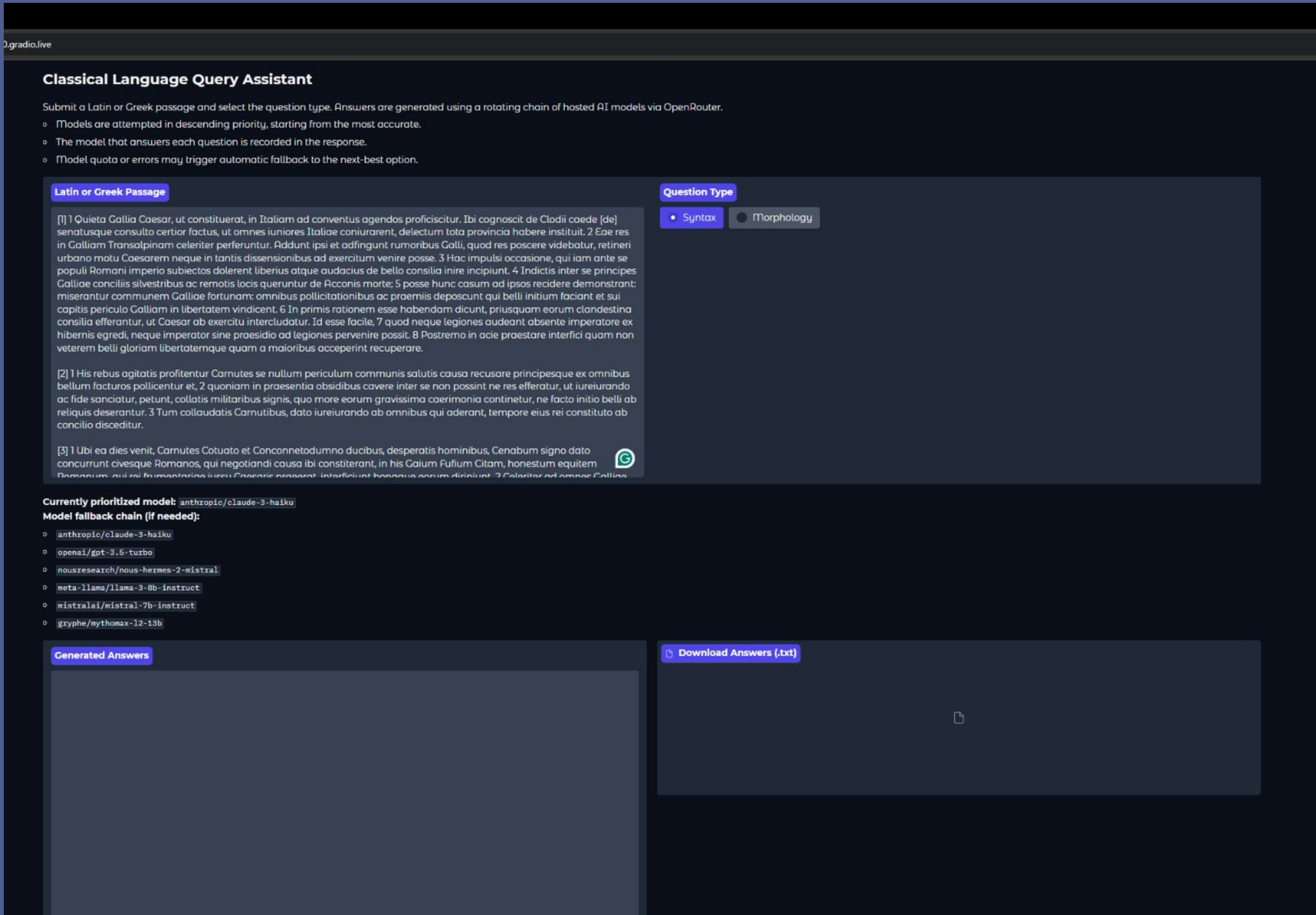


**Figure 3:** Query system interface

## Future Directions

Future directions include the further development of both the Greek and Latin lemmatizers, sentence modeling, and sentiment analysis models.

Furthermore, we are exploring the development of our own chatbot designed to interact with users in Ancient Greek and Latin. This initiative aims to bridge the gap between modern AI technologies and classical studies, fostering engagement with ancient languages in innovative ways.

**\*Data Driven Humanities Team Members:**

Justine Asman, Wavid Bowman, Eleni Bozia, Aidan Burrowes, Thomas Cerniglia, Srija Dey, Gebril Fradj, Nicole Fong, Guhan Gnanam, Zach Hracho, Jacob Hoppenstedt, Jeffrey James, Eden Layman, Jonnhy Liu, Daniel Miller, Connor Munjed, Niketha Nethaji, Krish Patel, Gillian Rodgers, Shae Robinson, Trey Slaten, Abraham Stefanos, Jordan Yu