

InterIDEAS: Philosophical Intertextuality via LLMs

| | | |
|---|--|--|
| Yue Yang* Maincode Monash University yue@maincode.com | Yinzhi Xu* Nanyang Technological University The University of Chicago yinzhi.xu@ntu.edu.sg | Chenghao Huang* Monash University chenghao.huang@monash.edu |
|---|--|--|

| | | |
|--|--|--|
| JohnMichael Jurgensen The University of Chicago jmjurg@uchicago.edu | Han Hu Independent Researcher agmaiofhuan@gmail.com | Hao Wang Monash University hao.wang2@monash.edu |
|--|--|--|

Abstract

The formation and circulation of ideas in philosophy have profound implications for understanding philosophical dynamism—enabling us to identify seminal texts, delineate intellectual traditions, and track changing conventions in the act of philosophizing. However, traditional analyses of these issues often depend on manual reading and subjective interpretation, constrained by human cognitive limits. We introduce InterIDEAS, a pioneering dataset designed to bridge philosophy, literary studies, and natural language processing (NLP). By merging theories of intertextuality from literary studies with bibliometric techniques and recent LLMs, InterIDEAS facilitates both quantitative and qualitative analysis of the intellectual, social, and historical relations in authentic philosophical texts. This dataset not only assists the study of philosophy but also contributes to the development of language models by providing a training corpus that enhances their interpretative capacity. The code URL for the dataset is https://github.com/interIDEAS/InterIDEAS_data.

1 Introduction

Although philosophy seems to be produced independently by a few genius thinkers, ideas do not exist in a vacuum. Philosophers read, cite, and discuss each other. Intertextuality—the relationship among different texts established by their referencing to or commenting on each other—is one of the most crucial ways to situate an idea in its epistemological, disciplinary, and social contexts. An adequate interpretation of even a single philosophical concept requires the reading of a vast collection of texts to understand with whom the philosopher(s) conversed, what sociohistorical incidents they responded to, and what intellectual foundation they evoked.

Previous researchers have addressed intertextuality via bibliometrics (Hammarfelt, 2016; Glänzel and Schoepflin, 1999): quantitatively analyzing citation entries, scholars can measure the relationships among texts and gain broad insights about a topic or even an entire discipline. However, directly extracting bibliographies from philosophy texts is not feasible, unless we limit ourselves to a very specific domain and to texts produced in a narrow span of time (Ahlgren et al., 2015). First, the lack of standardized citation practices before the mid-twentieth century results in a wide variety of formats that automated systems struggle to interpret. Second, the density of philosophical writing imposes tremendous challenges for digitalizing and processing.

A typical intertextual case in philosophy may read as follows: “The striving toward phenomenology was present already in the wonderfully profound Cartesian fundamental considerations; then, again, in the psychologism of the Lockean school; Hume almost set foot upon its domain, but with blinded eyes. And then the first to correctly see it was Kant, whose greatest intuitions become wholly understandable to us only when we had obtained by hard work a fully clear awareness of the peculiarity of the province belonging to phenomenology” (Husserl and Moran, 2012, p.142). Many factors contribute to the obscurity of this passage: a series of names, references, and concepts are crammed into a narrow space; the author writes rhetorically; the author does not specify his opinion to each mentioned philosopher and expects readers to uncover logical connections throughout the passage based on their previous philosophical knowledge; moreover, seemingly unimportant words like “almost” and “only” radically alter the author’s attitude. All this subtlety needs to be addressed, organized, and analyzed through a specifically designed data extraction process in order to organically integrate data-driven approaches into philosophical research.

*contributed equally to this work.

To address these challenges, we propose a novel data collection approach to curate a comprehensive dataset called InterIDEAS. We will show its workflow and structure that integrate LLMs’ reading capacity and human expertise. Venturing beyond usual bibliometric techniques that only analyze well-formulated citation entries, our prompt schema structures authentic philosophical writings in a manner that is organizable and analyzable by LLMs without effacing philosophers’ subtle reasoning. We systematically evaluate the RAG framework for complex information extraction and labeling tasks in philosophy, showing that LLM-based pipelines can substantially reduce costly expert annotation while achieving accuracy levels comparable to human experts.

Besides various domain-specific applications that we will demonstrate in this paper, our work contributes conceptual insights to interdisciplinary studies. The successful application of LLMs to philosophy, a discipline that values originality, individual voices, and subjective judgements, implies that dichotomies like qualitative and quantitative analyses, personal genius and general intellectual trends, textual details and immense data, as well as intimate reading and machinery processing of texts do not go against each other. They work in tandem to reveal lacunae overlooked by traditional methodologies in both the humanities and the sciences.

2 Related Works

Inquiry in intertextuality has been manually conducted by sociologists of philosophy like Randall Collins, who plotted network diagrams depicting philosophers’ personal relationships, educational affiliations, and intellectual lineages according to his own extensive reading (Collins, 2009). However, the innately limited recollection, speed, and processing of human reading subject Collins’ project to criticism like bias in text selection and interpretative methodologies.

Research in other disciplines provides novel avenues to address these issues. On the quantitative side, gathering and cross-comparing bibliographies in scientific and social scientific writings, bibliometrics offers ways to track relations among texts and achieve panoramic insights. For instance, given a specific topic and time period, we can ask how the frequency of its discussions changes over time, which articles are considered central or marginal, and the like (Leydesdorff and Ams-

terdamska, 1990). On the qualitative side, even though there is not a consensus regarding literary scholars’ taxonomy for references, there are plenty of concepts enabling us to describe the semantic structure, rhetorical impact, and implications of each reference with subtlety (Hohl Trillini and Quassdorf, 2010).

Humanities scholars have employed data-driven approaches and natural language processing (NLP) in studying dense writing, investigating topics like patterns in titles (Moretti, 2009) and abstracts (Ahlgren et al., 2015), evolution of a field (Bonino et al., 2022), authorial attribution (Peng and Hengartner, 2002), computational representation of arguments (Thagard, 2018), etc. A few pioneering datasets in intertextuality for humanities fields include *Hyperhamlet* (a database gathering a corpus of references to Hamlet in literature (Hohl Trillini and Quassdorf, 2010)), *Digital Dante* (a database mapping relations among writings by Dante and Ovid (Van Peteghem, 2020)), and *EDHIPHY* (a database extracting Anglo-American philosophers’ mentioning of each other in academic publications (Petrovich et al., 2024)). However, in the first two examples, relations are drawn from a few texts to address very specific research interests. In the third case, while mentions are vital for macroscopic relational networks and indexical purposes, they cannot support more qualitative analysis; for the database only record the frequency of mentions, effacing their content and purposes.

The employment of traditional transformers explains the specificity and even narrowness of the projects mentioned above. Traditional transformers’ limitations in restricted context understanding, poor reasoning capabilities, and limited knowledge integration forced researchers to confine themselves to textual details with limited, unambiguous markers. Recent advances in LLM such as GPT-3, T0, Galactica and LLaMa have demonstrated significant developments in NLP (Sanh et al., 2021; Touvron et al., 2023; Taylor et al., 2022); moreover, GPT-4 possesses notably enhanced capabilities in contextual understanding and reasoning. These abilities have been leveraged to manufacture textual datasets. For instance, the NORMDIAL dataset explores social norm adherence and violations in dialogue systems, using LLMs to generate culturally contextual conversations and thus pushing the boundaries of cross-cultural language modeling (Li et al., 2023). PoemSum (Mahbub et al., 2023) tests

LLMs’ ability to summarize poetry while retaining deeper figurative meanings.

Although LLMs have proven effective in dataset manufacturing and other NLP tasks in a diverse range of settings (Chang et al., 2024; Hu et al., 2023a,b, 2024), their application in niche humanities areas, such as philosophy, is less examined. Thus, in this work, we propose a framework that integrates prompt tuning, retrieval-augmented generation (RAG), and HITL examination to generate answers for intertextuality-related questions based on philosophical texts. Our dataset approaches intertextuality through semantic interpretation of authentic philosophical texts, moving beyond making comparisons at the word level and gathering statistics according to predetermined keywords and already formulated content. LLMs’ effective processing of texts and their generative nature enable us to devise a descriptive and evaluative schema, collect copious references including their content, function, and attitude reflected in detailed word and syntax choices, and envision the dataset’s applicability in both philosophy and AI.

3 Cross-Referential Data Collection

Our goal is to enable the selected LLM, namely GPT-4, to capture patterns and handle cross-referential data in philosophical texts via RAG and prompt engineering. The data include references—ranging from casual mentions and quotations to extensive critiques—of people, cultural productions, historical events, and social groups in modern philosophy. This section introduces our workflow for teaching the LLM to extract textual details with precision while avoiding hallucination. We first use RAG to convert texts into representations that aid contextual understanding. Next, prompt engineering guides LLMs to generate more accurate answers, while experts iteratively refine the prompts based on feedback. To understand how each step enhances performance, we include an ablation study in Appendix J. We validate the effectiveness of the framework by evaluating the quality of the resulting data set.

3.1 Data Collection Workflow

Fig. 1 illustrates the workflow of our data collection process. The vector base functions as the retrieval module in RAG (Lewis et al., 2020), enabling the LLM to access external knowledge during text generation. The process begins with a philosopher’s

book: 1) The text is segmented into chunks ①, embedded into vectors via a text encoder ②, and stored in a vector base ③; 2) When querying, relevant vectors are retrieved ④, combined with engineered QA prompts (White et al., 2023) for enhanced effectiveness, and passed to the LLM to extract reference attributes (detailed in Section 4.1) ⑤; 3) Philosophy experts review and analyze LLM outputs to iteratively refine prompts ⑥. Final high-quality QA pairs are stored in the database ⑦.

3.1.1 Philosophical Text Processing

To standardize input, all texts are first converted into PDF and split into semantically coherent paragraphs to fit the LLM’s context window while preserving local reference context. Each reference is described using three parameters—content type, intertextual function, and sentiment. These parameters are selected because of their relevance to philosophical research and their LLM-evaluability (see Limitations and Appendix L for details). On top of this standardized representation, we adopt a Retrieval-Augmented Generation (RAG) pipeline, in which the segmented texts are individually embedded into dense vector representations using the text-embedding-ada-002 model provided by OpenAI. This model outputs 1536-dimensional vectors widely applied in semantic similarity tasks. The resulting embeddings are stored and indexed in Chroma, an open-source vector database optimized for similarity search. During inference, queries are embedded using the same model, and the GPT API (version gpt-4-0314) retrieves the top-k relevant segments based on cosine similarity, thereby grounding generation in contextually aligned evidence.

3.1.2 Prompt Engineering

To enhance response quality, we employ three techniques (Fig. 2) as shown below. More prompt examples are provided in Appendix I:

- Role-Playing (RP): The LLM assumes the role of a philosopher (Static Info in Fig. 2), generating expert-style answers.
- Chain of Thought (CoT) (Wei et al., 2022): Questions are decomposed into sequential reasoning steps—starting with identifying references (upper blue box in Fig. 2), followed by evaluating the three attributes (lower boxes in Fig. 2).

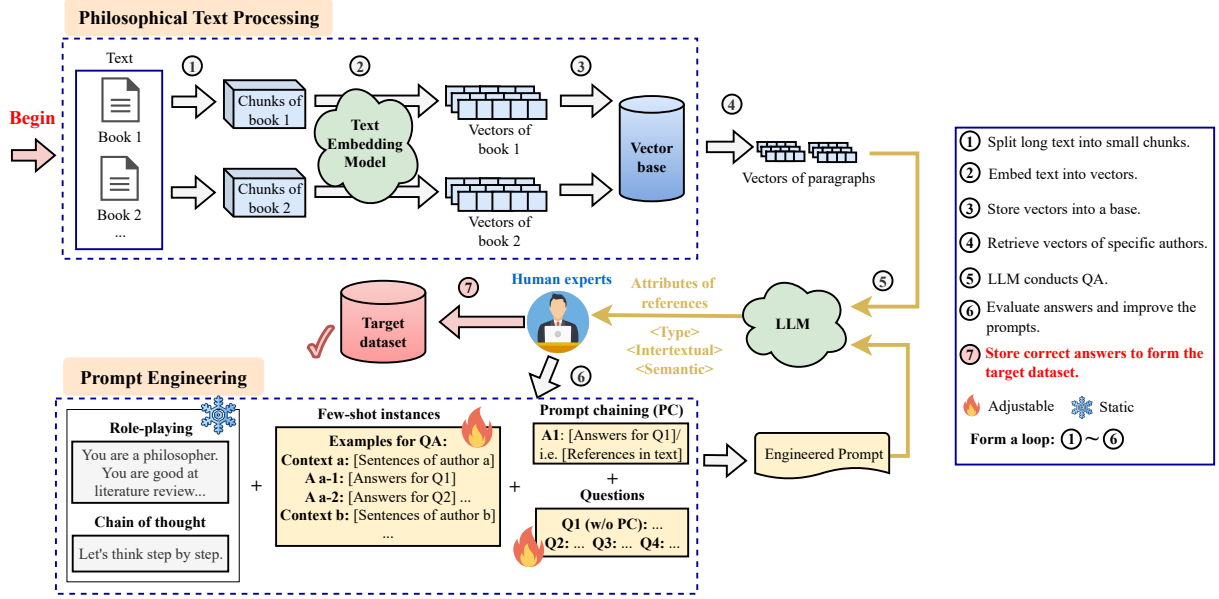


Figure 1: The entire workflow of the proposed data collection framework.

- Few-Shot Prompting (FS) (Brown et al., 2020): Contextual examples and corresponding answers (Few-shot Instances and Answer of Instances in Fig. 2) guide the model in interpreting the task.

3.1.3 Answer Evaluation and Prompt Improvement by Human Expert

To address limitations of the LLM, a dedicated expert review phase iteratively refines prompts and corrects recurrent mistakes. Experts assess LLM responses, identify common failure patterns, and incorporate them into prompts as constraints or illustrative few-shot cases when necessary. Each time the LLM provides answers to a set of texts, human experts evaluate their accuracy and identify patterns in the errors. These identified patterns are then integrated into the respective question prompt as additional conditions. When the identified patterns of errors are difficult to express within a few words, the sentences will be added to few-shot instances as representative cases.

3.2 Data Quality Evaluation

To confirm the accuracy and showcase the efficacy of our approach in facilitating the comprehension of philosophical texts, we assess and contrast the proficiency of our approach with that of human experts, humanities students, non-humanities students, and LLM-only approaches in extracting references from materials (approximately 500 words each) sourced from modern philosophy. All these

are excerpts from canonical texts carefully curated due to the richness and complexity of their intertextual references.

In our experiment, human experts are individuals who have obtained advanced degrees in fields such as literature or philosophy. The group of students with bachelor’s degrees in the humanities (BoH in Table 1) consists of individuals who have and only have obtained a bachelor’s degree in fields like literature or philosophy. The other student cohort includes native and non-native English speakers attending college to study the sciences, possessing a wide range of English language proficiency levels. For the purpose of this study, we recruited 5 human experts, 16 humanities students and 29 students of other backgrounds in both Australia and the United States, aiming to ensure a diverse and representative sample of participants for a comprehensive comparison of information extraction capabilities across different demographic groups. LLM-only approaches include GPT3.5, GPT3.5 with few-shot examples, GPT4 and GPT4 with few-shot examples. At the outset of the experiment, all participants received comprehensive instructions outlining the experimental requirements. They were then tasked with identifying and categorizing all references within a given paragraph in a strict timeframe of 20 minutes. Performance is measured by recall and accuracy, then compared with human results. We adopt a common scale: $\text{Recall} = \frac{x}{y}$ and $\text{Accuracy} = \frac{x}{r}$, whereas x is the correct answers found, y is the answers given, and

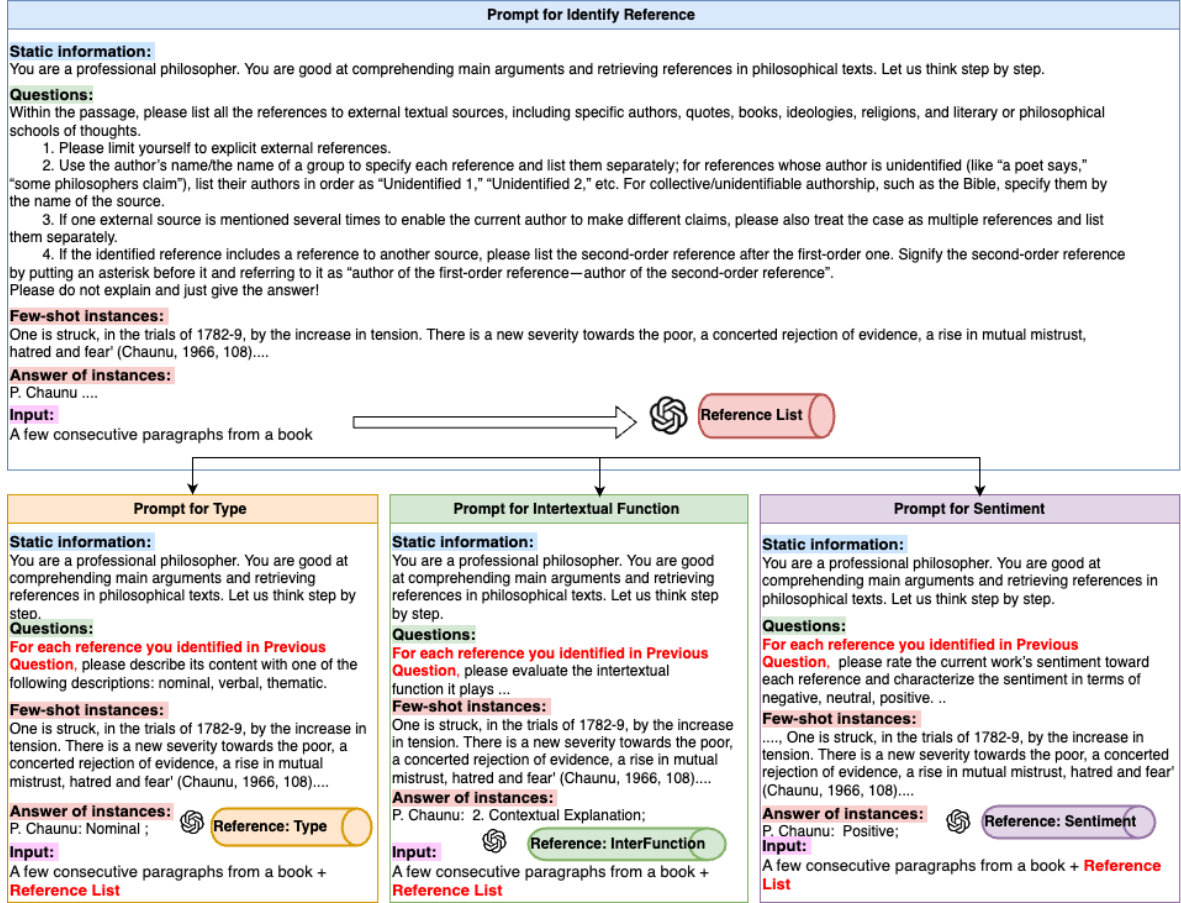


Figure 2: Prompting the LLM through few-shot examples to identify references, and evaluate their types, intertextual functions, and sentiments.

r is the total correct answers.

Table 1 shows the experimental results. Rows labeled Student/w.BoH, GPT3.5/w.FS, and GPT4/w.FS in the table correspond to the experimental results for the baselines: students with a Bachelor’s degree in Humanities, GPT-3.5 using few-shot examples, and GPT-4 using few-shot examples, respectively. Columns P_1 through P_6 in the table detail the accuracy and recall results for all baselines and our method, as applied to experiments on philosophical materials 1 through 6. Human experts outperform others, with amateurs struggling to grasp complex texts. Our approach ranks just below experts, excelling in accuracy and recall measures the model’s correct responses, indicating its precision. Although human experts achieve superior extraction outcomes compared to our method, the resource of human experts is extremely limited and costly. Thus, the experimental results verify that our method is effective, efficient, and economic, particularly in processing large-scale philosophical texts.

4 InterIDEAS Dataset Overview

In this dataset, we focus on philosophical books between 1750 and 1950. All were originally written in or have been translated into English. To date, we have analyzed over 45,000 pages. Still expanding, our dataset has amassed over 15,000 cross-referential data pairs, encompassing more than 3,150 philosophers and philosophical schools. Our periodization corresponds to the so-called “modern period” in the humanities—the period loosely bound by the beginning of the Industrial Revolution (circa 1760) and the end of WWII (1945). We slightly extended the timeline to address the time lag between historical events and their intellectual stimuli and reactions. In selecting texts, we balanced coverage with representativeness. We incorporated authors and texts into the dataset according to three objectives: 1) Covering prominent thinkers; 2) Featuring different geographical locations for intellectual debates, including traditional cultural centers like France, emerging intellectual hubs at that time like the U.S., and marginalized places like

Table 1: Evaluation matrix. Bold numbers indicate the highest results from P_1 - P_6 following human experts.

| | Accuracy | | | | | | Recall | | | | | |
|----------------|-------------|-------------|------------|-------------|-------------|-------------|-------------|-------------|------------|-------------|-------------|-------------|
| | P_1 | P_2 | P_3 | P_4 | P_5 | P_6 | P_1 | P_2 | P_3 | P_4 | P_5 | P_6 |
| Human Experts | 1 | 1 | 1 | 0.92 | 0.89 | 0.98 | 1 | 1 | 1 | 1 | 0.93 | 1 |
| Student/w.BoH | 0.97 | 0.75 | 0.63 | 0.75 | 0.75 | 0.64 | 0.85 | 0.74 | 0.79 | 0.66 | 0.56 | 0.71 |
| Other Students | 0.75 | 0.6 | 0.68 | 0.47 | 0.44 | 0.75 | 0.69 | 0.62 | 0.68 | 0.47 | 0.25 | 0.60 |
| GPT3.5 | 0.46 | 0.58 | 0.66 | 0.71 | 0.67 | 0.63 | 0.54 | 0.61 | 0.53 | 0.47 | 0.25 | 0.43 |
| GPT3.5/w.FS | 0.75 | 0.55 | 0.71 | 0.63 | 0.8 | 0.75 | 0.69 | 0.55 | 0.53 | 0.41 | 0.50 | 0.60 |
| GPT4/w.FS | 0.75 | 0.64 | 0.6 | 0.65 | 0.83 | 0.74 | 0.69 | 0.64 | 0.6 | 0.77 | 0.63 | 0.66 |
| Ours | 0.85 | 0.91 | 0.8 | 0.74 | 0.75 | 0.84 | 0.85 | 0.91 | 0.8 | 0.81 | 0.75 | 0.88 |

India; 3) Presenting works from authors of different occupations, including academics, journalists, political activists, novelists, and literary critics.

4.1 Metadata Format

Empirically speaking, most discussions of external materials in philosophy fall into the following categories: ideas or activities of specific agents or groups. Therefore, we delineate intertextuality as references to other discourses, including cultural productions (including all items involving intellectual efforts like books, artistic works, ideologies, and religions) and historical events (including both famous monumental events in history and words and deeds of other people in general). With our deliberately loose definition guiding the LLM to extract references of diverse nature—ranging from published texts to anecdotes, from specific individuals to vague social groups—the dataset reflects different philosophical, political, historical, and personal components that jointly contribute to the vibrancy of modern philosophy.

Table 2: Distribution of sentiments across intertextual functions.

| Intertextual Function | Negative | Neutral | Positive |
|-------------------------------------|----------|---------|----------|
| Name-dropping | 514 | 6537 | 778 |
| Contextual Explanation | 284 | 2626 | 657 |
| Critical Engagement | 620 | 2361 | 394 |
| Conceptual Application or Expansion | 12 | 119 | 145 |

Table 3: Distribution of sentiment types across reference categories.

| Type/Sentiment | Nominal | Thematic | Verbal |
|----------------|---------|----------|--------|
| Negative | 927 | 369 | 134 |
| Neutral | 7923 | 2713 | 1013 |
| Positive | 1376 | 420 | 184 |

Our dataset presents the bibliographic entries

for all the processed books. Each *Book Title* is connected to a group of *Reference Names*. Linked directly to each reference is the *Content Type* capturing its content and level of specificity. The content type assumes one of the following forms: the nominal, the verbal, or the thematic. “The nominal” identifies names of cultural productions, people, social groups, and events mentioned in the reference; “the verbal” records direct quotations from other texts; “the thematic” provides brief summaries for loose, unspecified discussion of external references. Each reference is associated with an *Intertextual Function*, which describes the rhetorical intent of the reference ranging from “name-dropping” (ND) and “contextual explanation” (CE_x) to “critical engagement” (CE_n) and “conceptual application or expansion” (CAoE). This classification helps us understand the extent of interaction between the current work and the referred content. Furthermore, the *Sentiment* attribute assesses the current author’s attitude towards each reference, which is categorized as “negative,” “neutral,” or “positive.” The relationships among these values are structured to ensure an one-to-one correspondence between a reference and its content type, intertextual function, and sentiment. The metadata schema can be found in Appendix L.

Based on our dataset, nominal references are the most common, constituting 67.9% of the data, followed by thematic and verbal references. In sentiment analysis, neutral sentiments predominate at 77.4%, with positive and negative sentiments at 13.1% and 9.5% respectively. For intertextual functions, name-dropping is most frequent, making up 52% of the instances, whereas critical engagement and contextual explanation are also significant, and conceptual application or expansion is relatively rare. Relevant visualizations can be found in Appendix K. These statistics illustrate the dominance of nominal referencing and neutral

sentiments in the dataset, with name-dropping being the primary intertextual function. Meanwhile, authors’ attitudes are crucial in determining the depth of their engagement with others’ ideas and actions, as shown in Table 2 and Table 3: negative attitudes are often suggested by explicit criticism; people, events, or works about which the current authors feel impartial tend to be cursorily discussed. The general statistics of our dataset also uncover features of modern philosophical writings. First, the dominance of neutral and positive sentiments shows that the field is largely organized by amicability. Second, the distribution of sentiments across intertextual functions suggests that in constructing philosophical arguments, philosophers generally adapt the style of discussion (“function” in the dataset) rather than the choice of materials (“type”) to reflect their attitudes (“sentiment”). Comparing the number of positive references with that of the negative ones, we find that philosophers express amicability more overtly and more frequently.

5 Applications of InterIDEAS in Philosophy and LLMs

5.1 InterIDEAS for Philosophy

Automating reference extraction lets philosophy researchers visualize how ideas propagate across centuries, schools, and authors—something infeasible by hand at scale. In this manner, we can reveal synchronic and diachronic patterns in philosophy. As a demonstration, we extract the 50 most frequent references that appear in at least 3 texts processed by our model. The word map 3a confirms the interdisciplinary nature of philosophy. Besides acclaimed philosophers and philosophical schools, we find religions (e.g., “Christianity,” “God,” “Buddha,” and “The Bible”) and political events and entities (e.g., “Roman Empire,” “British Empire,” and “French Revolution”) constitutive to philosophical discussion. We extract all individual philosophers in these common references. Most of them were active in the Mediterranean region and the English Channel region, belonging to one of the following three intellectual periods: ancient, enlightenment, and modern philosophy. We map a network Fig. 3b to visualize this flow of ideas. The network shows, for example, how likely a modern philosopher who has referred to Solon would also be influenced by Voltaire and moreover, by Schopenhauer. Our chart suggests that two important intellectual traditions for modern philosophy

are Plato-Rousseau-Hegel and Plato-John Stuart Mill-Engels.

In addition to general intellectual environment, the dataset assists our analysis of individual philosophers. For instance, by statistically presenting the proportion of a few selected authors’ attitudes in Fig. 4, we identify possible similarities in the tones of their writings: Georg Jellinek and Franz Oppenheimer may share a more placid style, while Russell’s writing tends to be more polemical. Further zooming into the intertextual network of these philosophers, we may discover previously unknown relationships. As shown by Fig. 3c, while Bertrand Russell and Émile Faguet are rarely discussed together, their shared strong sentiment for Homer and against John Stuart Mill casts light on their comparability. It further proposes possible incompatibility between Homer and Mill, due to which the commitment to one’s stance entails the rejection of the other’s.

5.2 InterIDEAS for LLMs

While our dataset is motivated by key concerns in philosophy, it contributes to NLP research. We will use sentiment analysis to demonstrate this potential: on the one hand, sentiment analysis has pervasive applications, which will allow us to join fundamental conversations happening in NLP; on the other hand, as a standard NLP benchmark, sentiment analysis attests to the quality of our dataset through its potential in model fine-tuning. We construct 2,236 reference–attitude pairs suitable for sentiment classification. Each pair comprises a sentence from an authentic philosophical text and its author’s assessed attitudes towards the referenced content. These pairs are divided into training (70%), validation (20%), and test sets (10%), where in the test set, samples with label “negative”, “neutral”, and “positive” are 142, 53, and 33. We consider not only LLMs but also pre-trained language models (PLMs) in our experiment. PLMs focus on pre-training to generate general language representations for downstream tasks, while LLMs primarily focus on natural language generation and typically involve larger model scales. Since both models can be fine-tuned to adapt to downstream tasks, we select five popular PLMs and four outstanding LLMs for fine-tuning. The five PLMs can be categorized into three types: 1) BERT-based: BERT (Devlin et al., 2018), ALBERT (Lan et al., 2019), and BERTweet (Nguyen et al., 2020); 2) RoBERTa (Liu et al., 2019); 3) XLNet (Yang et al.,

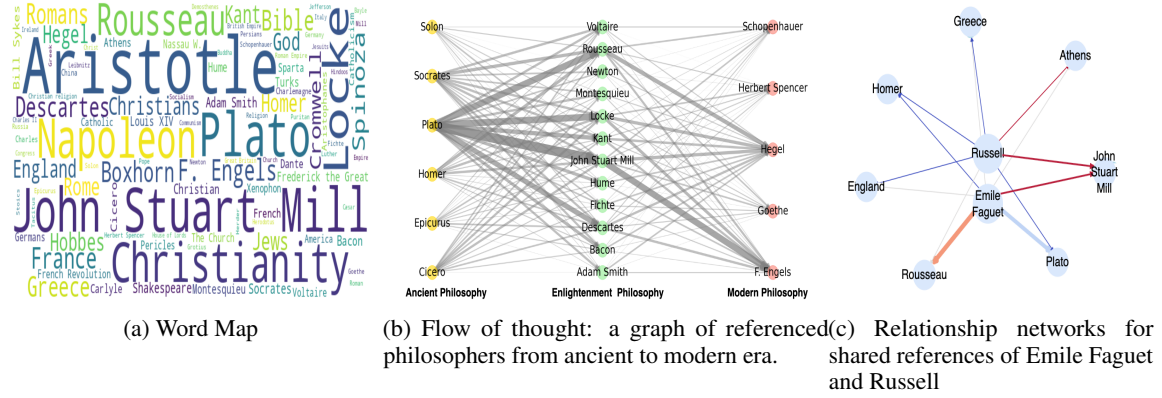


Figure 3: Philosophical references analysis

Table 4: Popular open-source PLMs and LLMs for sentiment classification on the proposed dataset w./w.o. fine-tuning, or few-shot learning for GPT-4.

| Model | Before fine-tuning/few-shot | | | | After fine-tuning/few-shot | | | | Computational cost | | |
|----------------|-----------------------------|--------------|--------------|--------------|----------------------------|--------------|--------------|--------------|--------------------|-------|------|
| | Acc. | F1 | Pre. | Rec. | Acc. | F1 | Pre. | Rec. | Param. | FT % | Sec. |
| BERT | 16.67 | 14.24 | 28.36 | 30.26 | 63.32 | 39.01 | 51.59 | 39.69 | 0.11B | 1.21% | 69 |
| ALBERT | 14.91 | 9.72 | 16.00 | 33.96 | 60.96 | 25.25 | 20.59 | 32.63 | 0.05B | 0.24% | 32 |
| BERTweet | 28.51 | 22.28 | 36.44 | 34.94 | 60.96 | 34.23 | 37.48 | 36.57 | 0.13B | 0.98% | 61 |
| RoBERTa | 23.25 | 12.57 | 7.75 | 33.33 | 63.16 | 45.68 | 50.80 | 44.76 | 0.12B | 2.00% | 222 |
| XLNet | 28.07 | 24.73 | 37.81 | 38.19 | 49.56 | 35.48 | 35.45 | 35.54 | 0.12B | 0.62% | 245 |
| Average | 22.28 | 16.67 | 25.27 | 34.14 | 59.59 | 35.93 | 39.18 | 37.84 | - | - | - |
| Llama 2 | 26.75 | 25.39 | 35.52 | 29.17 | 62.28 | 53.17 | 54.03 | 52.49 | 6.54B | 0.50% | 677 |
| Llama 3 | 27.63 | 27.79 | 40.82 | 39.77 | 67.54 | 62.61 | 61.02 | 65.45 | 7.51B | 0.52% | 747 |
| Mistral | 25.88 | 25.59 | 32.68 | 36.81 | 50.44 | 45.20 | 45.30 | 49.98 | 7.11B | 0.94% | 859 |
| GPT-2 | 27.19 | 27.72 | 41.90 | 39.08 | 53.95 | 48.42 | 47.41 | 51.11 | 0.38B | 0.88% | 175 |
| Average | 26.86 | 26.62 | 37.73 | 36.21 | 58.55 | 52.35 | 51.94 | 54.76 | - | - | - |
| GPT-4 | 24.56 | 21.03 | 34.91 | 33.58 | 42.54 | 40.79 | 51.05 | 47.47 | - | - | - |

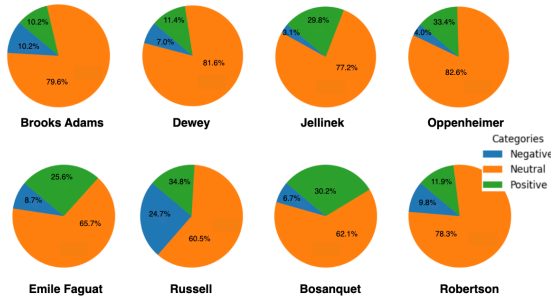


Figure 4: Sentiment pie chart for selected authors

2019). The four LLMs can be classified into three types: 1) Llama-based: Llama 2-7B and Llama 3-8B (Touvron et al., 2023); 2) Mistral-7B (Jiang et al., 2023); 3) GPT-2 (Radford et al., 2019). Additionally, we study GPT-4o (Achiam et al., 2023), which is the most state-of-the-art (SOTA) LLM, to do direct inference without any extra training. We randomly choose 5 samples of each label from the training set as the few-shot instances for GPT-

4o. We track how the performance of all PLMs and LLMs pre-trained for text/sequence classification changes, before and after fine-tuning on our reference-attitude dataset for 100 epochs.

Evaluation metrics include accuracy, macro F1 score, macro precision, and macro recall, to calculate more reasonable results of the imbalanced test set. Additionally, the size of each model, the proportion of fine-tuned parameters, and the time cost for fine-tuning are recorded in Table 4. The confusion matrices of each model are shown and analyzed in Appendix M.

In Table 4, the average performance improvements after fine-tuning are noteworthy. The average accuracy of PLMs and LLMs increased from 22.28% and 26.86% to 59.59% and 58.55%, and the average F1 score improved from 16.67% and 26.62% to 35.93% and 52.35%, respectively. This demonstrates that our provided philosophical corpus exhibits significant potential for fine-tuning. Among all these models, Llama 3 achieves the

best performance, which can be attributed to its large number of parameters. Overall, the accuracy of PLMs is generally slightly higher than that of LLMs, but the F1 scores are noticeably lower. This could be attributed to the fact that PLMs have significantly fewer parameters than LLMs, coupled with the presence of data imbalance in the training set (with more negative samples). As a result, overfitting during fine-tuning PLMs might have occurred, causing the outputs to be heavily biased towards the negative class. PLMs consume less computational resources compared to LLMs. This indicates that PLMs, while less resource-intensive, may struggle with achieving balanced performance across different classes in the context of imbalanced datasets, particularly in complex tasks like sentiment analysis of philosophical texts. Additionally, the results from GPT-4 show that even simple few-shot learning markedly improves output quality. This validates the representational quality of our dataset samples.

We present confusion matrices of Llama 3 w./w.o. fine-tuning, GPT-4o w./w.o. few-shot learning, and Mistral w./w.o. fine-tuning in Fig. 5. Before fine-tuning or few-shot learning, all three models tend to favor a single class: Llama 3 and Mistral lean toward the positive class, while GPT-4o is initially biased toward the neutral. Notably, none consistently predicts the negative class, despite its abundance in the test set. This phenomenon highlights that even advanced LLMs suffer from prediction unbalance when directly applied to philosophical sentiment classification. After fine-tuning, both Llama 3 and Mistral exhibit a marked shift toward the negative class, indicating that fine-tuning improves alignment with the dominant distribution but sometimes at the expense of neutral and positive recognition. GPT-4o, however, demonstrates the most stable and balanced performance: although it initially favors the neutral, with few-shot learning it improves across all three classes rather than collapsing into one. This suggests that our philosophical corpus may be especially effective in few-shot settings, where models like GPT-4o can leverage its representational richness without the drawbacks of overfitting.

6 Conclusion

In this article, we introduce InterIDEAS for extracting and analyzing philosophical intertextuality. Enhanced by both LLMs and human expertise,

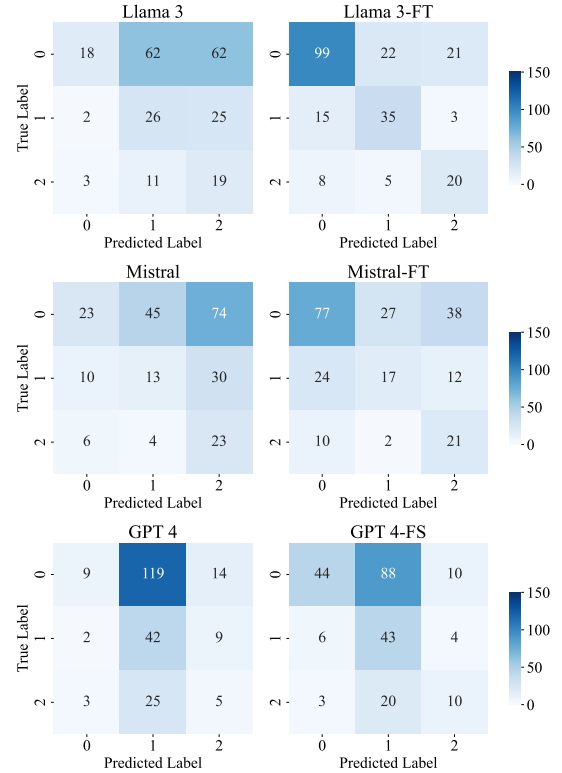


Figure 5: Confusion matrices of Llama 3 w./w.o. fine-tuning and GPT-4o w./w.o. few-shot learning adopted for sentiment classification.

InterIDEAS provides a robust foundation for exploring intellectual structures and dynamics. This dataset inspires philosophical scholars to venture beyond traditional research methodologies, uncover unknown relations among texts and writers, and understand philosophy not only as a set of texts and ideas but also as a vital field of interacting forces driven by personal innovations, interpersonal connections, disciplinary conventions, and cultural conditions. Our work further advances NLP and AI by tackling the unique challenges of processing texts with abstract meaning, ambiguous syntax, lengthy reasoning processes, and complex logic that exceed conventional benchmarks.

In the future, we aim to extend the dataset to encompass a wider range of intellectual activities, languages, and historical periods, further investigating how our approach and dataset can be employed to study other text-based, intertextually-dense fields like law and history.

Limitations

Intertextuality can be implicit and unnoticed even by its users. This project restricts itself to direct references with clear textual evidence for two reasons. First, a collection of all probable intertextuality will lead to the postmodernist belief that everything is intertextual. A profound idea conceptually speaking, it might be an aimless data collecting method that will result in distracting, low-quality data. Second, while our data suggest novel directions for philosophical research, the in-depth verification and interpretation of these preliminary insights will often require researchers to return to at least some of the original textual evidence.

We organize intertextual instances into three fixed dimensions to ensure that their most important features are captured, i.e., content, function, and sentiment. Yet we are aware that these dimensions cannot exhaust all possible rhetorical sophistication of intertextuality. We employ this arguably artificial structure to make the collected data comparable to each other, so that valuable patterns will not be obscured by a nebula of particularities.

Limitations of using LLMs for processing philosophical texts found in our work are summarized as follows: 1) *Semantic dissection*: When multiple references are listed in parallel grammatical structures, the LLM may categorize them into different functions, even though they assume identical rhetorical roles. Through manual review, representative sentences are integrated into few-shot instances, and some constraints are imposed on the questions, effectively mitigating this issue. 2) *Literal-mindedness*: The LLM struggles in literary expressions with complex emotions, such as rhetorical questions and irony. This aspect has seen some improvement through the addition of few-shot instances. 3) *Stereotyping*: Faced with specific input information, such as “Hitler,” the LLM tends to respond based on its built-in stereotypes with “negative” disregarding the author’s potentially “neutral” or “positive” stance.

Limitations of our dataset include: 1) *Style*: The dataset excludes symbol- and aphorism-based texts, which require the designing of a completely different approach to parse, collect, and analyze their intertextuality. Since symbols tend to be heavily featured in philosophical subfields like logic and philosophy of language, and since certain philosophers like Wittgenstein have a predilection for aphorisms, our dataset can potentially exclude a few topics and

writers. 2) *Language*: Our current approach is limited to texts that are written in or have been translated into English. This limitation can raise concerns of Eurocentrism. To address these problems, we hope to extend the approach to other styles and languages in the future by recruiting philosophical researchers with different research and language expertise.

Acknowledgments

We express our deep gratitude to Professor Haun Saussy at the University of Chicago for his indispensable intellectual and financial support, when we first embarked on this project. We gratefully acknowledge Maincode (Australia) for their generous support, which significantly contributed to both the development and dissemination of this work.

Yue Yang and Yinzhi Xu jointly designed the approach, conducted the experiments and the analysis, and drafted the main text. Chenghao Huang made a significant contribution to the pipeline of the dataset. JohnMichael Jurgensen provided important insights at the project’s early stage and proofread this paper. Han Hu helped with editorial revisions and structural improvements of this paper. Hao Wang offered valuable guidance throughout the project.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Per Ahlgren, Peter Pagin, Olle Persson, and Maria Svedberg. 2015. Bibliometric analysis of two subdomains in philosophy: Free will and sorites. *Scientometrics*, 103:47–73.
- Guido Bonino, Paolo Maffezoli, Eugenio Petrovich, and Paolo Tripodi. 2022. When philosophy (of science) meets formal methods: a citation analysis of early approaches between research fields. *Synthese*, 200(2):177.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, and 1 others. 2024. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45.
- Randall Collins. 2009. *The sociology of philosophies*. Harvard University Press.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Wolfgang Glänzel and Urs Schoepflin. 1999. A bibliometric study of reference literature in the sciences and social sciences. *Information processing & management*, 35(1):31–44.
- Björn Hammarfelt. 2016. Beyond coverage: Toward a bibliometrics for the humanities. *Research assessment in the humanities: Towards criteria and procedures*, pages 115–131.
- Regula Hohl Trillini and Sixta Quassdorf. 2010. A ‘key to all quotations’? a corpus-based parameter model of intertextuality. *Literary and Linguistic Computing*, 25(3):269–286.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Han Hu, Ruiqi Dong, John Grundy, Thai Minh Nguyen, Huaxiao Liu, and Chunyang Chen. 2023a. Automated mapping of adaptive app guis from phones to tvs. *ACM Transactions on Software Engineering and Methodology*, 33(2):1–31.
- Han Hu, Han Wang, Ruiqi Dong, Xiao Chen, and Chunyang Chen. 2024. Enhancing gui exploration coverage of android apps with deep link-integrated monkey. *ACM Transactions on Software Engineering and Methodology*, 33(6):1–31.
- Han Hu, Haolan Zhan, Yujin Huang, and Di Liu. 2023b. Pairwise gui dataset construction between android phones and tablets. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pages 59860–59872.
- Edmund Husserl and Dermot Moran. 2012. *Ideas: General introduction to pure phenomenology*. Routledge.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, and 1 others. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Loet Leydesdorff and Olga Amsterdamska. 1990. Dimensions of citation analysis. *Science, Technology, & Human Values*, 15(3):305–335.
- Oliver Li, Mallika Subramanian, Arkadiy Saakyan, Sky CH-Wang, and Smaranda Muresan. 2023. [NormDial: A comparable bilingual synthetic dialog dataset for modeling social norm adherence and violation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15732–15744, Singapore. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Ridwan Mahbub, Ifrad Khan, Samiha Anuva, Md Shihab Shahriar, Md Tahmid Rahman Laskar, and Sabbir Ahmed. 2023. [Unveiling the essence of poetry: Introducing a comprehensive dataset and benchmark for poem summarization](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14878–14886, Singapore. Association for Computational Linguistics.
- Franco Moretti. 2009. Style, inc. reflections on seven thousand titles (british novels, 1740–1850). *Critical Inquiry*, 36(1):134–158.

- Dat Quoc Nguyen, Thanh Vu, and Anh Tuan Nguyen. 2020. Bertweet: A pre-trained language model for english tweets. *arXiv preprint arXiv:2005.10200*.
- Roger D Peng and Nicolas W Hengartner. 2002. Quantitative analysis of literary styles. *The American Statistician*, 56(3):175–185.
- Eugenio Petrovich, Sander Verhaegh, Gregor Bös, Claudia Cristalli, Fons Dewulf, Ties van Gemert, and Nina IJdens. 2024. Bibliometrics beyond citations: introducing mention extraction and analysis. *Scientometrics*, 129(9):5731–5768.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and 1 others. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, and 1 others. 2021. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*.
- Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*.
- Paul Thagard. 2018. Computational models in science and philosophy. *Introduction to formal philosophy*, pages 457–467.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Julie Van Peteghem. 2020. Ovid in dante’s commedia. In *Italian Readers of Ovid from the Origins to Petrarch*, pages 169–222. Brill.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C Schmidt. 2023. A prompt pattern catalog to enhance prompt engineering with chatgpt. *arXiv preprint arXiv:2302.11382*.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.

A Licensing

All the data we currently open to public are originating from Project Gutenberg <https://gutenberg.org/about/>. Project Gutenberg eBooks may be freely used in the United States because most are not protected by U.S. copyright law. They may not be free of copyright in other countries. Readers outside of the United States must check the copyright terms of their countries before accessing, downloading or redistributing eBooks. We also have a number of copyrighted titles, for which the copyright holder has given permission for unlimited non-commercial worldwide use. For Project Gutenberg, no permission is needed for non-commercial use. So, for example, you can freely redistribute any eBook, anywhere, any time, with or without the “Project Gutenberg” trademark included. The “Small Print” has more details. Note that if you are not in the US, you must confirm yourself whether an item is free to redistribute where you are.

The copyright status of philosophy books can vary significantly depending on several factors, such as the date of the author’s death and the specific laws of the country in which the book was published. Here are some general guidelines: In most countries, works enter the public domain 70 years after the death of the author. If the author of a philosophy book died more than 70 years ago, it is likely that their works are now in the public domain. Besides, some philosophy books, especially classic texts, may be in the public domain, but newer editions (which might include modern commentary, translations, or annotations) can still be protected by copyright. Copyright laws can vary from one country to another. For example, some countries have extensions for certain types of works or authors.

For the remaining unpublished data, we are actively working on verifying the copyright status and obtaining the necessary permissions. We will continue to update our dataset as soon as we confirm the copyright status of each book and secure the appropriate permissions.

B Accuracy for Whole Dataset

Given the lack of available tools other than human expertise for verifying the accuracy of the resulting dataset, and considering the impracticality of human experts reviewing all responses due to the extensive volume of material, we have adopted a strategy of randomly selecting 5 text chunks per

100 for manual verification. Additionally, we plan to make this dataset accessible for future research use and will provide an interface allowing users to identify errors and update the dataset accordingly. Based on the random sample and check, GPT showed remarkable precision in recognizing 98.11% of references to external sources across all books. Additionally, it was able to accurately depict 93% of the content from these identified references. As of the current date, language learning models (LLMs) have achieved a 75.7% success rate in identifying intertextual functions and an 86.4% success rate in sentiment analysis.

At this stage, our goal is to confirm that the performance of the LLM is stable across texts. Verifying its performance on a random 5% pages for each book we processed is sufficient to reflect its overall performance. Meanwhile, 5% of 45000 pages is 2250 pages. Each of our human experts spent on average 10 minutes reading a page, processing 15-20 pages per day. 5% is already a taxing workload.

C Human Reading Capability Experiment

C.1 Instructions

Objective: The aim of this experiment is to assess the intertextual reading ability of individuals at various levels of proficiency. Participants will be asked to read texts of differing complexity and respond to the listed questions. We focus on assessing LLM performance against general human performance, not just versus experts. We include both expert and non-expert readers of philosophical texts. The results show that LLMs perform better than nonprofessionals, though they fall short of expert levels. This suggests that our dataset can expand experts' analytic scope and improve nonprofessionals' understanding of textual details. It also implies that the task requires specialized knowledge or skills that are beyond the capacity of general participants and highlights the effectiveness of the LLM in handling complex scenarios where typical human capabilities are insufficient. Such findings might be essential for understanding the limits of human performance in specific contexts and the potential areas where advanced models like LLMs can be particularly beneficial.

Participant Requirements:

- Age: 20-80
- Language Proficiency: Participants must be

college students or individuals with higher education, residing in an English-speaking country.

Materials Provided

- A series of texts at varying levels of difficulty.
- A questionnaire for each text to assess intertextual reading ability.

C.1.1 Procedure

Introduction: Participants will receive an overview of the experiment, including its purpose and what will be required of them.

Consent: Participants must read and sign a consent form agreeing to partake in the experiment and acknowledging the confidentiality and use of their data.

Pre-Test Survey: A short survey to gather participant background information relevant to the study, such as age, education level, and reading habits.

Pre-Reading: Participants will give 15 minutes to read the instruction for questions

Reading Task: Participants will be given one or two texts. Each text should be read in a quiet environment without distractions. Participants are advised to read at their natural pace.

Comprehension Assessment: After reading each text, participants will answer a set of questions. The questions may be multiple choice, short answer, or a mix of both.

Breaks: Participants are allowed to take short breaks between texts if needed.

Post-Reading Survey: After completing all the readings, participants will fill out a survey capturing their experience, challenges faced, and any feedback on the texts.

Debriefing: Participants will be provided with a summary of the experiment and its objectives. Any questions or concerns from participants will be addressed.

C.1.2 Ethics and Confidentiality

All participant information will be kept confidential. Participants have the right to withdraw from the study at any point without any negative consequences.

C.1.3 Contact Information

Provide contact details for participants to reach out if they have any questions or concerns before, during, or after the experiment. Thank you for your participation and valuable contribution to this research!

C.1.4 Compensation

Each participant is provided with a \$15 coupon for the school coffee shop.

C.2 Questions

C.2.1 Q1 for Reference Identification

Within the passage, please list all the references to external textual sources, including specific authors, quotes, books, ideologies, religions, and literary or philosophical schools of thoughts. Use the author's name/the name of a group to specify each reference; for references whose author is unidentified (like "a poet says," "some philosophers claim"), list their authors in order as "Unidentified 1," "Unidentified 2," etc. For collective/unidentifiable authorship, such as the Bible, specify them by the name of the source.

C.2.2 Q2 for Content Type

For each reference you identified, please describe its content with one or more of the following descriptions: 1. Nominal, meaning those references that explicitly mention names of other authors, books, collections of works, and other schools of thought in the main text; for nominal references, signal their content by exact names used in the passage. If there are multiple nominal references, separate them by colons. E.g., Marx: nominal (Marx; The Communist Manifesto) 2. Verbal, meaning direct quotation of phrases and sentences from other sources; for verbal references, signal their content by abbreviated versions of the quotes that only keep the first and the last two words of the quote, with ellipses in between. If there are multiple verbal references, separate them by colons. E.g., Marx: verbal ("the history... class struggles") 3. Thematic, meaning references to others' claims, ideas, and motifs not through direct quotes but through paraphrases; for thematic references, please signify their content by a summary in one or two philosophical terms. If there are multiple thematic references, separate them by colons. E.g., Marx: thematic (child labor)

C.2.3 Q3 for Intertextual Function

For each reference identified in prompt 1, please evaluate the intertextual function it plays by the closest descriptions below. Classify the references by "Name-Dropping," "Contextual Explanation," "Critical Engagement," or "Conceptual Application or Expansion." 1. Name-Dropping: This category is for when the current work merely mentions

the names of authors, works, or concepts as representative cases of a phenomenon or an argument, without detailed explanations. 2. Contextual Explanation: Elements of external sources are mentioned and given some exposition to clarify the source's relevance to the author's argument. These references add depth to the discussion but are presented without the author's personal judgment of the reference as right or wrong. Examples include references to factual evidence in support of the argument, references that intend to exemplify the author's arguments, etc. 3. Critical Engagement: In this category, the current work actively engages with external sources by offering detailed analysis (at least one sentence of analysis for each reference) and value judgements. The author's subjective attitudes are evident as they express their agreements or disagreements with the ideas presented in the reference. 4. Conceptual Application or Expansion: References that fall into this category are not only explained but are also used as a springboard for further development of the current work.

C.2.4 Q4 for Sentiment

Please rate the current author's sentiment toward each reference identified in prompt 1, and characterize the sentiment in terms of strongly negative, negative, neutral, positive, strongly positive. If the author's attitude is ambiguous or unknown, please label it as "neutral". For references to historical facts, please label them as "neutral". Organize your final answer as: Marx Nominal (Marx; The Communist Manifesto); Verbal ("the history... class struggles"); Thematic (child labor) 3. Critical Engagement Positive

C.3 Ethical Approval

This study was approved by the Monash University Human Research Ethics Committee (MUHREC) in accordance with the National Statement on Ethical Conduct in Human Research (2023) and the principles of the Declaration of Helsinki. The approved project was granted ethics approval under Project ID: 44944 (Review Reference: 2024-44944-115031).

D Static Analysis for the Data Quality Evaluation

D.1 Accuracy

Human Experts have the highest consistency with an average score of 0.965 and a standard deviation of 0.044. Their performance distribution may

Table 5: Summary of accuracy results with statistical analysis.

| Group | Scores | | | | | | Average | Std. Dev. | P-value |
|----------------|--------|------|------|------|------|------|---------|-----------|---------|
| Human Experts | 1 | 1 | 1 | 0.92 | 0.89 | 0.98 | 0.965 | 0.044 | 0.039 |
| Student/w.BoH | 0.97 | 0.75 | 0.63 | 0.75 | 0.75 | 0.64 | 0.748 | 0.112 | 0.110 |
| Other Students | 0.75 | 0.6 | 0.68 | 0.47 | 0.44 | 0.75 | 0.615 | 0.124 | 0.258 |
| GPT3.5 | 0.46 | 0.58 | 0.66 | 0.71 | 0.67 | 0.63 | 0.618 | 0.081 | 0.382 |
| GPT3.5/w.FS | 0.75 | 0.55 | 0.71 | 0.63 | 0.8 | 0.75 | 0.698 | 0.084 | 0.523 |
| GPT4/w.FS | 0.75 | 0.64 | 0.6 | 0.65 | 0.83 | 0.74 | 0.702 | 0.079 | 0.659 |
| Ours | 0.85 | 0.91 | 0.8 | 0.74 | 0.75 | 0.84 | 0.815 | 0.059 | 0.722 |

not be normal (p-value = 0.039). Student with BoH shows moderate variability with an average of 0.748 and a standard deviation of 0.112, with performance deemed normally distributed (p-value = 0.110). Other Students have the most variability with an average of 0.615 and a standard deviation of 0.124, and normal distribution (p-value = 0.258). GPT3.5 and GPT3.5 with FS score averages of 0.618 and 0.698, respectively, both with normal performance distributions (p-values > 0.380). GPT4 with FS and GPT4 with FPEh show consistent high performance with averages of 0.702 and 0.815, respectively, and low variability (SD < 0.08), with normal distribution (p-values > 0.650).

D.2 Recall

The updated dataset table presents a comprehensive statistical analysis of performance scores from various groups, including Human Experts, Students with and without Book of Humanities (BoH), and different versions of GPT models. The Human Experts group exhibits nearly perfect scores with an average of 0.988 and a minimal standard deviation of 0.026, although their scores do not follow a normal distribution. In contrast, the Student groups show more variability, with averages of 0.718 and 0.552 for Students with BoH and Other Students, respectively. The GPT models display a progression in performance from GPT3.5 to our approach with GPT4, where the latter achieves an impressive average of 0.833 with a standard deviation of 0.053, showing a more consistent performance (normality p-value = 0.955).

E Interview with Human Experts

We further surveyed human experts about their opinions on our dataset. All of our human experts, who are either university professors of philosophy or PhD students in the humanities, find this dataset

both intriguing and valuable. Representing a bridge between traditional academic studies and the latest technological advancements, our application offers a novel method for integrating these two fields. One of our interviewees said, “Given the vast scope of work that no individual could complete in a lifetime, the use of language learning models now makes this formidable task feasible.” Another interviewee recognized the philosophical implication of our approach: “Philosophy is a strange field, with a style of inquiry sometimes behaving like mathematics and sometimes like literary studies. The seeming incompatibility between the two sets of assumptions is what keeps me coming back to it, and this investigation clarifies a lot.” One professor was intrigued by how our approach gives concrete guidance for practical pedagogical tasks like designing syllabus and creating analytical assignments by showing the interrelations among texts. A PhD student pointed out that the granularity of the information in the dataset is “just right”; the dataset provides crucial clues to interpretation and further learning, without reductive summaries that may discourage students from reading the actual texts.

F Data Format for Fine-Tuning

To illustrate the utility of the proposed dataset in natural language processing and data science, a sentiment classification dataset containing 2,236 entries has been developed. Each entry includes a sentence from philosophical texts, accompanied by the author’s expressed sentiment towards the referenced content within that sentence, as follows 6:

G Computational Resources

All data collection processes and fine-tuning experiments are conducted on a server with 8 NVIDIA

Table 6: Summary of recall results with statistical analysis.

| Group | Scores | | | | | | Average | Std. Dev. | P-value |
|----------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|--------------|------------------|
| Human Experts | 1 | 1 | 1 | 1 | 0.93 | 1 | 0.988 | 0.026 | $2.07 * 10^{-5}$ |
| Student/w.BoH | 0.85 | 0.74 | 0.79 | 0.66 | 0.56 | 0.71 | 0.718 | 0.093 | 0.985 |
| Other Students | 0.69 | 0.62 | 0.68 | 0.47 | 0.25 | 0.60 | 0.552 | 0.153 | 0.135 |
| GPT3.5 | 0.54 | 0.61 | 0.53 | 0.47 | 0.25 | 0.43 | 0.472 | 0.114 | 0.487 |
| GPT3.5/w.FS | 0.69 | 0.55 | 0.53 | 0.41 | 0.50 | 0.60 | 0.547 | 0.086 | 0.987 |
| GPT4/w.FS | 0.69 | 0.64 | 0.60 | 0.77 | 0.63 | 0.66 | 0.665 | 0.054 | 0.518 |
| Ours | 0.85 | 0.91 | 0.80 | 0.81 | 0.75 | 0.88 | 0.833 | 0.053 | 0.955 |

Table 7: Hyperparameters details.

| Module | Parameter | Parameter description | Value |
|-------------|------------------------|-------------------------------|---|
| LoRA | r_{LoRA} | The rank of LoRA matrix | 8 |
| | α_{LoRA} | Scaling factor of LoRA matrix | 32 |
| | δ_{LoRA} | Dropout rate | 0.1 |
| | | | If XLNet: [layer_1, layer_2] |
| | | | elif Llama or Mistral: |
| Fine-tuning | θ_{LoRA} | Modules to be fine-tuned | [q_proj, k_proj, v_proj, o_proj, gate_proj, up_proj, down_proj] |
| | | | elif GPT-2: [c_attn, c_fc, c_proj] |
| | | | else: [query, key, value, dense] |
| | r | Learning rate | 1e-4 |
| | E | Training epoch | 100 |
| | γ | Weight decay | 0.01 |
| | B | Batch size | 16 |

```

Instruction: Please rate the current work sentiment toward 'Montessori system' and characterize the sentiment in terms of
negative, neutral, positive
context: These educational theorists who have
had a knowledge of children, such as the inventors of Kindergarten and
the Montessori system, [14] have not always had enough realization of
the ultimate goal of education to be able to deal successfully with
advanced instruction.
response: Neutral
category: closed_q

```

Figure 6: Data format for fine-tuning.

GeForce 3090 GPUs, each of which has 24G memory. The CUDA version is 11.5.

All the resource usage for sentiment classification through fine-tuning is presented in Table 4, including the model parameter count, the proportion of fine-tuned parameters to the total parameter count, and the time required for 100 epochs of fine-tuning. For details on the fine-tuning parameters, please refer to Table 7.

H Training details for Sentiment Classification

The sentiment classification fine-tuning runs based on Transformer package under Python 3.9, where

the version of Pytorch is 1.12. All models are downloaded from Huggingface, pre-trained on sentiment or emotion corpus ¹.

Data split: The dataset is split into training set (70%), validation set (20%), and test set (10%) with the random seed 42 and shuffling. Specially,

¹BERT: <https://huggingface.co/google-bert/bert-base-uncased>;
ALBERT: <https://huggingface.co/tals/albert-xlarge-vitaminc-mnli>;
BERTweet: <https://huggingface.co/cardiffnlp/bertweet-base-sentiment>;
RoBERTa: <https://huggingface.co/cardiffnlp/twitter-roberta-base-sentiment>;
XLNet: <https://huggingface.co/TehranNLP/xlnet-base-cased-mnli>;
Llama 2: <https://huggingface.co/Mikael110/llama-2-7b-guanaco-fp16>;
Llama 3: <https://huggingface.co/RLHFlow/ArmoRM-Llama3-8B-v0.1>;
Mistral: <https://huggingface.co/weqweasdas/RM-Mistral-7B>;
GPT-2: <https://huggingface.co/michelecafagna26/gpt2-medium-finetuned-sst2-sentiment>.

for BERTweet, the maximal length of each input sample is truncated to 128 due to the fixed model input dimension.

Hyperparameters: To reduce the computational cost of LLM fine-tuning, we adopt Low-Rank Adaptation (LoRA) (Hu et al., 2021) by Parameter Efficient Fine-Tuning (PEFT) package. For fine-tuning, we adopt Transformer Package. Both hyperparameters of LoRA and fine-tuning keep the same for all experimented models, recorded in Table 7. The hyperparameters corresponding to each model follow the default settings on Huggingface.

The rank r_{LoRA} is set to 8, determining the rank of the low-rank matrices used by LoRA. It affects the reduction in model parameters and computational efficiency by defining the dimension of the introduced low-rank matrices. The scaling factor α_{LoRA} is set to 32, controlling the scaling size of the adaptation matrices during training. By adjusting this factor, the magnitude of the adaptation matrices' updates can be balanced to avoid excessively large or small updates. The dropout rate δ_{LoRA} is set to 0.1, meaning that 10% of the neurons will be randomly dropped during training, helping prevent overfitting and enhances the generalization capability of the model. Last but not least, the particular modules θ_{LoRA} are specified to be fine-tuned. These hyperparameters work together to optimize the application of LoRA in specific models and tasks, balancing computational cost and model performance.

In terms of fine-tuning, the learning rate r is set to $1e-4$, determining the magnitude of updates to the model parameters at each step. A smaller learning rate ensures that the model updates its parameters in small, precise steps, contributing to a stable and refined training process, reducing the risk of instability from large parameter changes. The training epoch E is set to 100 to avoid underfitting but might lead to overfitting. To help with it, the weight decay rate is set to 0.01 by reducing the size of the model weights at each update. The batch size B is set to 16 due to both the size of our proposed sentiment classification dataset and our hardware limitation. Additionally, the optimizer is ADAM, and the load accuracy is 32 bit for all models.

I Prompts for References

We show all the prompts in Fig. 7, 8, 9, and 10:

Prompt for Reference 1

Static information:

You are a professional philosopher. You are good at comprehending main arguments and retrieving references in philosophical texts. Let us think step by step.

Question:

Within the passage, please list all the references to external textual sources, including specific authors, quotes, books, ideologies, religions, and literary or philosophical schools of thoughts.

1. Please limit yourself to explicit external references.
2. Use the author's name/the name of a group to specify each reference and list them separately; for references whose author is unidentified (like "a poet says," "some philosophers claim"), list their authors in order as "Unidentified 1," "Unidentified 2," etc. For collective/unidentifiable authorship, such as the Bible, specify them by the name of the source.
3. If one external source is mentioned several times to enable the current author to make different claims, please also treat the case as multiple references and list them separately.
4. If the identified reference includes a reference to another source, please list the second-order reference after the first-order one. Signify the second-order reference by putting an asterisk before it and referring to it as "author of the first-order reference—author of the second-order reference". Please do not explain and just give the answer!

Few-shot instances:

Context:

One is struck, in the trials of 1782-9, by the increase in tension. There is a new severity towards the poor, a concerted rejection of evidence, a rise in mutual mistrust, hatred and fear' (Chaunu, 1966, 108).

...

Homage is paid to the 'great reformers' - Beccaria, Servan, Dupaty, Lacretelle, Duport, Pastoret, Target, Bergasse, the compilers of the Cahiers, or petitions, and the Constituent Assembly - for having imposed this leniency on a legal machinery and on 'classical' theoreticians who, at the end of the eighteenth century, were still rejecting it with well-formulated arguments.

...

What is this nationalist political theory about? ... This is opposed to imperialism, which seeks to bring peace and prosperity to the world by uniting mankind, as much as possible, under a single political regime. ... At that time, the struggle against Communism ended, and the minds of Western leaders became preoccupied with two great imperialist projects ...

Answers of instances:

P. Chaunu; Beccaria, Servan, Dupaty, Lacretelle, Duport, Pastoret, Target, Bergasse; Imperialism; Communism; ...

Figure 7. The engineered prompt for the 1st question for references.

Prompt for Reference 2

Static information:

You are a professional philosopher. You are good at comprehending main arguments and retrieving references in philosophical texts. Let us think step by step.

Question:

For each reference you identified in question 1, please describe its content with one or more of the following descriptions:

1. Nominal, meaning those references that explicitly mention names of other authors, books, collections of works, and other schools of thought in the main text; for nominal references, signal their content by exact names used in the passage. Specification of authors or sources in citational practice does not count as nominal. If there are multiple nominal references, separate them by colons.

2. Verbal, meaning direct quotation of phrases and sentences from other sources; for verbal references, signal their content by abbreviated versions of the quotes that only keep the first and the last two words of the quote, with ellipses in between. If there are multiple verbal references, separate them by colons.

3. Thematic, meaning references to others' claims, ideas, and motifs not through direct quotes but through paraphrases; for thematic references, please signify their content by a summary in one or two philosophical terms. If there are multiple thematic references, separate them by colons.

If there is no reference to others' claims in a category, please give NA.

If one external source is mentioned several times to enable the current author to make different claims, please also treat the case as multiple references and list them separately.

Lastly, formulate your answer in this way:

Referred item: nominal (content of the nominal references); verbal (content of the verbal references);

3. thematic (content of the thematic references)

Please do not explain and just give the answer!

Few-shot instances:

In these few shot examples, we covered all the cases. When you run the prompt, please choose the most applicable one for each reference. You don't need to identify all functions within a passage.

These are examples for your answer:

Context:

The same as the context in Fig. 7.

Answers of instances:

P. Chaunu: Nominal (P. Chaunu); Verbal ("a constant... for security"); Thematic (crime; economic pressure);

Beccaria, Servan, Dupaty, Lacrosette, Dupont, Pastoret, Target, Bergasse: Nominal (Beccaria, Servan, Dupaty, Lacrosette, Dupont, Pastoret, Target, Bergasse, Cahiers);

Imperialism: Thematic (Alternative to nationalism);

Communism: Thematic (the Cold War);

...

Prompt for Reference 3

Static information:

You are a professional philosopher. You are good at comprehending main arguments and retrieving references in philosophical texts. Let us think step by step.

Question:

For each reference identified in question 1, please evaluate the intertextual function it plays by the closest descriptions below. Classify the references by "Name-Dropping," "Contextual Explanation," "Critical Engagement," or "Conceptual Application or Expansion";

1. **Name-Dropping:** This category is for when the current work merely mentions the names of authors, works, or concepts as representative cases of a phenomenon or an argument, without detailed explanations that exceed one sentence. In particular, if there is a list of names whose individual significance is not discussed, please label them as "Name-Dropping." Other markers for this category include mentioning in passing like "c.f.," "for details, please see..." etc.

2. **Contextual Explanation:** Elements of external sources are mentioned and given some exposition to clarify the source's relevance to the author's argument. These references add depth to the discussion but are presented without the author's personal judgment of the reference as right or wrong. Examples include references to factual evidence in support of the argument, references that intend to exemplify the author's arguments, etc.

3. **Critical Engagement:** In this category, the current work actively engages with external sources by offering detailed analysis (at least one sentence of analysis for each reference) and value judgements. The author's subjective attitudes are evident as they express their agreements or disagreements with the ideas presented in these references.

4. **Conceptual Application or Expansion:** References that fall into this category are not only explained but are also used as a springboard for further development of the current work. The current work distills keywords or arguments from the reference and expands upon them, possibly transforming them or integrating them into a new framework. Examples include a problematic concept that is adjusted and employed in further discussion; a methodology from other sources is adopted by the current author, etc.

If one external source is mentioned several times to enable the current author to make different claims, please also treat the case as multiple references and list them separately.

Please do not explain and just give the answer!

Few-shot instances:

In these few shot examples, we covered all the cases. When you run the prompt, please choose the most applicable one for each reference. You don't need to identify all functions within a passage.

These are examples for your answer:

Figure 8. The engineered prompt for the 2nd question for references.

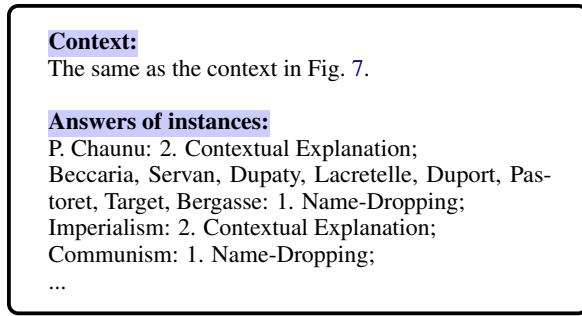


Figure 9. The engineered prompt for the 3rd question for references.

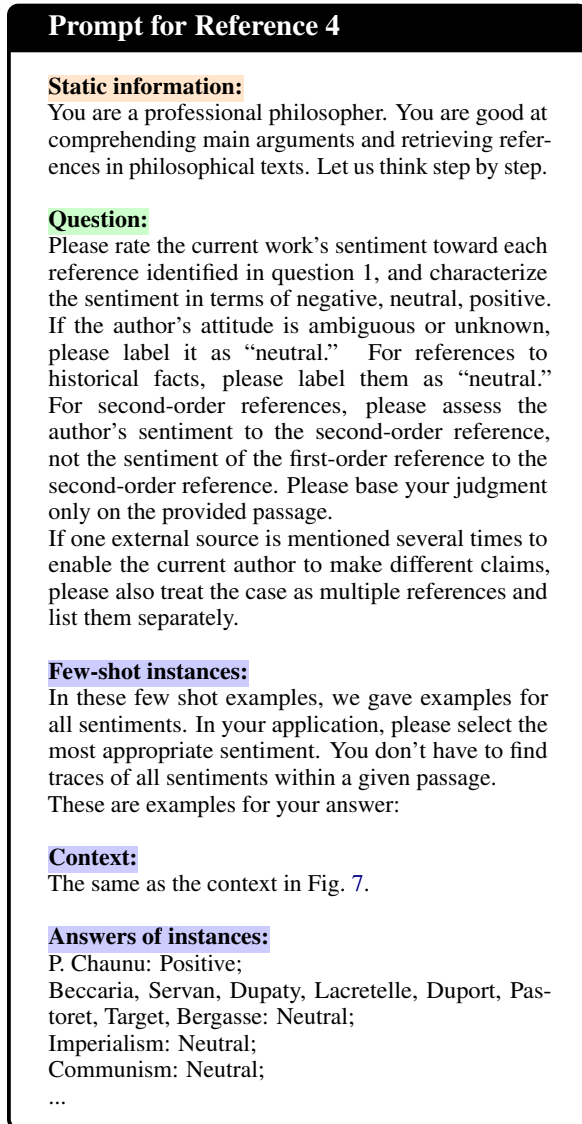


Figure 10. The engineered prompt for the 4th question for references.

J Ablation Study

To evaluate the contribution of each component in our framework, we conduct an ablation study comparing four experimental settings: (1) **GPT4 (baseline)**, (2) **GPT4 + RAG**, (3) **GPT4 + FS + RAG**, and (4) **Ours (GPT4 + FS + PE + RAG)**.

Performance is measured in terms of accuracy and recall across six paragraphs. The results are summarised in Figure 11.

The baseline (**GPT4**) achieves low accuracy (0.15–0.32) across paragraphs, reflecting its difficulty in extracting information without retrieval support. Introducing **RAG** produces a substantial improvement (up to 0.73), confirming that retrieval mechanisms significantly enhance factual grounding. Adding **FS** with RAG further stabilises results and boosts accuracy in most paragraphs (up to 0.83), suggesting FS helps reduce noise from irrelevant retrievals. Finally, our approach (**FS + PE + RAG**) achieves the highest accuracy overall (0.74–0.91), with prompt engineering synergising with FS and RAG to improve both peak performance and stability across paragraphs.

Baseline recall is weak (0.18–0.33), mirroring the accuracy problem. With **RAG**, recall improves significantly (up to 0.62), but performance varies across paragraphs. **FS + RAG** improves consistency, recovering dips (e.g., Paragraph 4 at 0.77). Our final model consistently achieves the best recall (0.75–0.91), showing that prompt engineering enhances recall in addition to accuracy by structuring how retrieved information is incorporated.

RAG provides the largest single improvement, establishing it as the critical foundation. FS contributes to stability and noise reduction, particularly in recall, while also improving upper-bound accuracy. Prompt engineering provides the decisive final lift, ensuring strong and reliable performance across all paragraphs. This ablation study highlights the progressive importance of RAG, FS, and PE. Without retrieval, GPT4 struggles to achieve meaningful accuracy or recall. With RAG, the model grounds its responses but suffers from variability. FS stabilises results, and prompt engineering pushes performance to the highest level, achieving state-of-the-art accuracy and recall across all paragraphs. The synergy of retrieval, feature selection, and prompt engineering is therefore essential for unlocking the full potential of large language models in complex information extraction tasks.

K Data Distribution

See distribution of the dataset in Fig. 12.

L Metadata format and description

We present the metadata schema specifically designed for the analysis of intertextual references

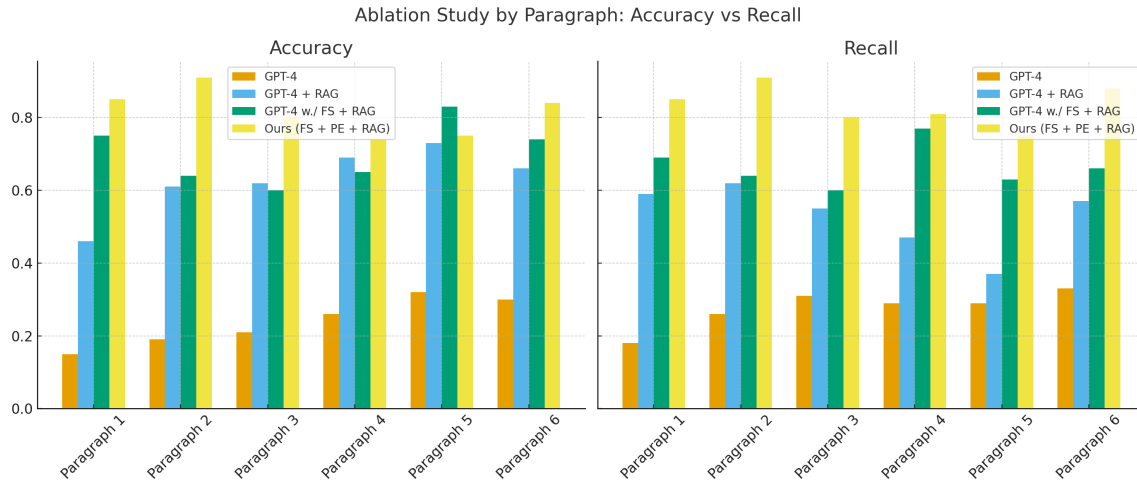


Figure 11: Ablation study results comparing accuracy (left) and recall (right) across six paragraphs.

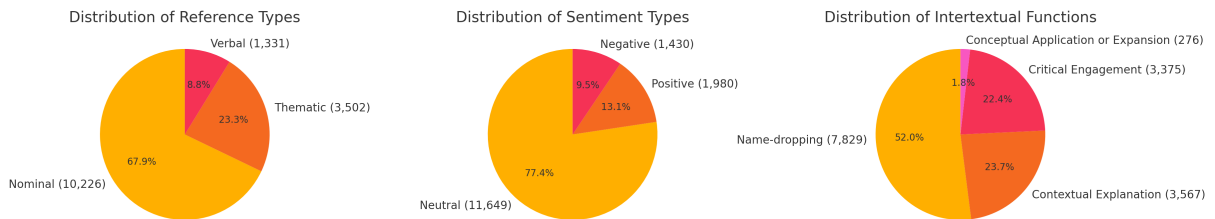


Figure 12: Pie charts showing the distribution of reference types, sentiment types, and intertextual functions.

within humanities writing (as mentioned in Section 4.1) in Fig. 13.

M Supplementary Analysis on Sentiment Classification

The confusion matrices of each PLM or LLM is shown in Figure 14. It can be observed that both PLMs and LLMs tend to output a specific class, as seen in the following patterns: Neutral - BERTweet, RoBERTa, XLNet, Llama 2, GPT-4; Positive - BERT, ALBERT, Llama 3, Mistral, GPT-2. Notably, none of the models consistently favors the Negative class, even though Negative samples are the most abundant in the test set. This tendency could be attributed to the differences in the pre-training corpora and methods used for each model. Additionally, LLMs exhibit more moderate biases compared to PLMs, especially in more recent models like Llama 3, which also has the largest number of parameters. This can be attributed to the enhanced language understanding capabilities of LLMs, driven by their larger parameter counts and more extensive training corpora. Nonetheless, this highlights a significant issue: even the most advanced language models suffer from severe predic-

tion unbalance when directly performing sentiment classification in a philosophical context. Therefore, the most straightforward approach to enhance a language model’s understanding of philosophical texts is fine-tuning.

After fine-tuning, it is evident that all models become more inclined to output Negative. To some extent, this suggests that the overall trend brought by fine-tuning is benefiting. However, this trend appears to be extreme, even impairing the models’ ability to correctly classify Neutral and Positive instances. This could be due to the imbalance in the training dataset. Similarly, the output bias in LLMs remains less pronounced than in PLMs, which can once again be attributed to the ability of LLMs to better handle imbalanced datasets due to their larger parameter counts.

GPT-4 demonstrates the most stable and balanced performance. Although GPT-4 initially leans towards Neutral, after few-shot learning, it shows improvement in predicting all three classes rather than favoring one. This may indicate that our corpus has greater potential when used for few-shot learning, perhaps even more so than for fine-tuning.

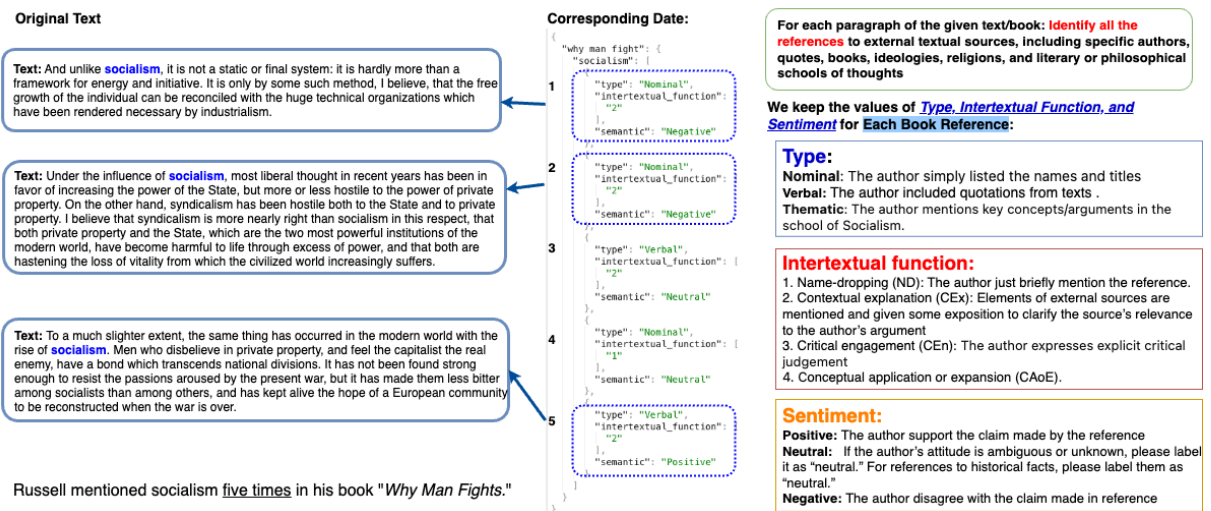


Figure 13: Metadata format and description.

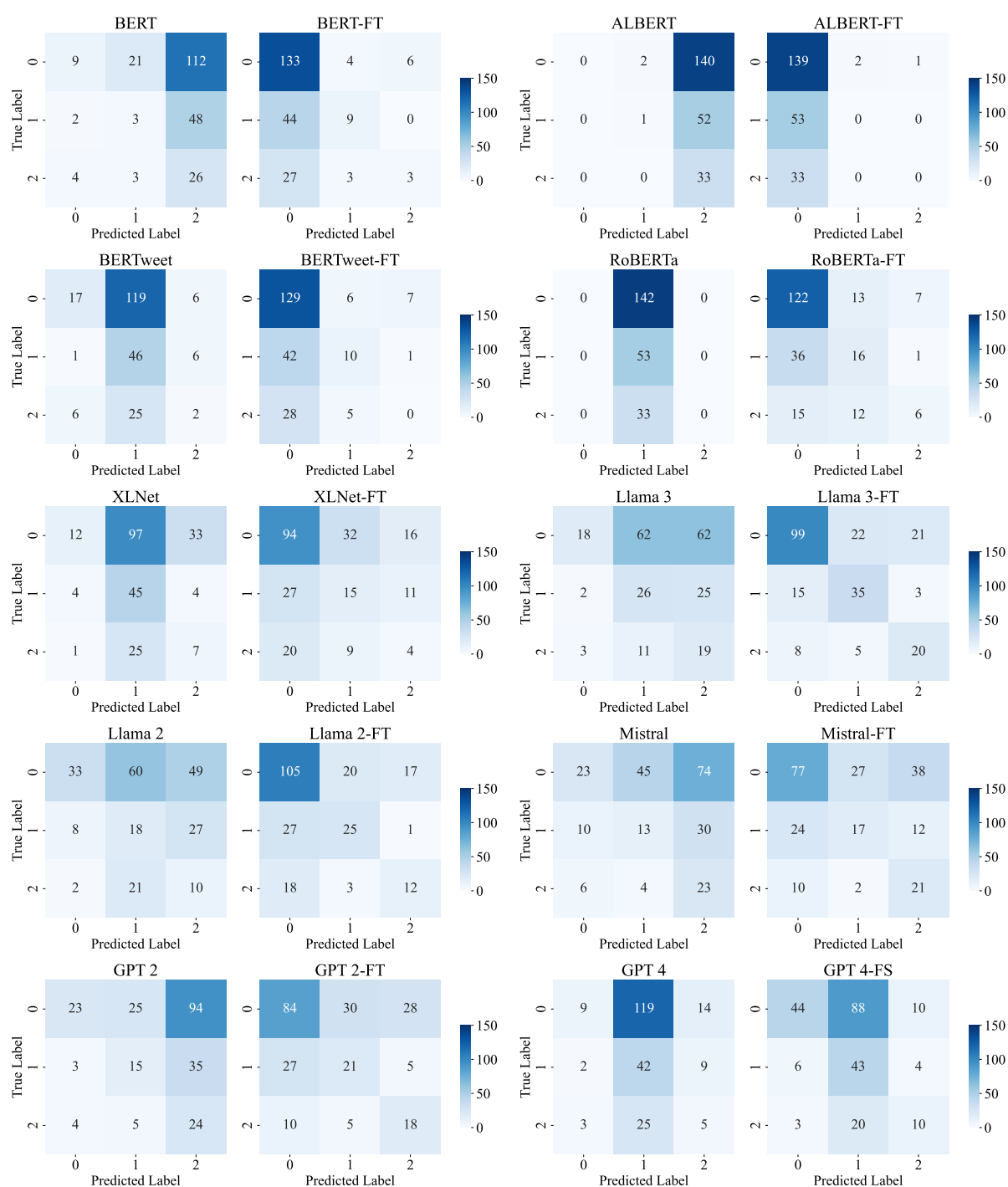


Figure 14: Confusion matrices of each model adopted for sentiment classification before and after fine-tuning or few-shot learning.