InterIDEAS: Philosophical Intertextuality via LLMs

Anonymous ACL submission

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 enables both quantitative and qualitative analysis of the intellectual, social, and historical relations embedded within 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 challenges and enhances their interpretative capacity.

1 Introduction

011

017

019

021

024

025

027

042

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 in philosophy, 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 digitalization and comparison. 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

For instance, 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.

To address these challenges, we propose a novel data collection approach to collect a comprehen-

084sive dataset called InterIDEAS. We will show its085workflow and structure, which integrates LLMs'086reading capacity and human expertise. Venturing087beyond usual bibliometric techniques that only an-088alyze well-formulated citation entries, our prompt089schema structures authentic philosophical writings090in a manner that is organizable and analyzable091by LLMs without effacing their subtle reasoning.092LLMs' successful application to philosophical in-093tertextuality will further imply their potential in as-094sisting research in other humanistic disciplines, like095literature and law, where the circulation and forma-096tion of ideas are encoded in stylistic language.

2 Related Works

097

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

125

126

127

128 129

130

131

132

133

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 measure relations among texts and achieve panoramic insights. For instance, given a specific topic and time period, we can investigate how the frequency of its discussions change over time, which articles are considered central or marginal, and the like (Leydesdorff and Amsterdamska, 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. 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

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

As demonstrated by the rather narrow scope or the specificity of some of the projects mentioned above, traditional transformers face limitations like restricted context understanding, poor reasoning capabilities, and limited knowledge integration-all of which create bottlenecks in humanities research that require deeper contextual analysis and crossdisciplinary insights. Recent advances in LLM such as GPT-3, T0, Galactica and LLaMa (Sanh et al., 2021; Touvron et al., 2023; Taylor et al., 2022) have marked significant developments in NLP, in which GPT-4, the latest product, has notably enhanced capabilities in language understanding, generation, and reasoning. These abilities have been leveraged to manufacture textual datasets that cast light on both humanities and AI research. For instance, the NORMDIAL dataset explores social norm adherence and violations in dialogue systems, using LLMs to generate culturally contextual conversations, pushing the boundaries of crosscultural language modeling (Li et al., 2023). Poem-Sum (Mahbub et al., 2023) tests LLMs' ability to summarize poetry while retaining deeper figurative meanings.

Although LLMs have proven effective in NLP dataset manufacturing and other general NLP tasks (Chang et al., 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 on philosophical texts. Our dataset approaches intertextuality through semantic interpretation of full texts of authentic philosophical writings, moving beyond making comparisons at the word level and gathering statistics according to predetermined

234

235

237

keywords and already formulated content. LLMs'
effective comprehension of texts and their generative nature enable us to devise a descriptive and
evaluative schema, to collect copious references
including their content, function, and attitude reflected in detailed word and syntax choices, and
to construct a dataset with flexible applicability in
both philosophy and AI.

3 Cross-Referential Data Collection

194

210

211

212

213

214

215

216

217

218

219

223

224

226

233

Our goal is to enable LLMs to capture patterns and 195 handle cross-referential data in philosophical texts via RAG and prompt engineering. The data in-197 cludes references-ranging from casual mentions 198 and quotations to extensive critiques-of people, 199 texts, and groups in political philosophy. This sec-200 tion introduces our workflow for teaching LLMs with philosophical texts while avoiding hallucination. We first use RAG to convert texts into representations that aid contextual understanding. Next, 204 prompt engineering guides LLMs to generate more accurate answers, while experts iteratively refine prompts based on feedback. Last but not least, we validate the framework's effectiveness by evaluating the quality of the resulting datasets.

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 LLMs 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 converted into PDF and split into paragraphs to fit the LLM's context window, preserving local reference context. Each reference is described using three parameters, including content type, intertextual function, and sentiment, selected for their argumentative relevance and LLM-evaluability.

3.1.2 Prompt Engineering

To enhance response quality, we employ three techniques (Fig. 2) 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).
- 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 LLM limitations, a dedicated expert review phase iteratively refines prompts and correct recurrent mistakes made by the LLM. 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, this study is designed to assess and contrast the proficiency of our collection approach with that of human experts, humanities students, other students and LLM-only approaches in identifying and extracting detailed information from philosophic materials (approximately 500 words each) sourced from modern philosophy.

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



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

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 ChatGPT3.5, Chat-GPT3.5 with few-shot examples, ChatGPT4 and ChatGPT4 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.

284

287

291

296

297

298

Performance is measured by recall and accuracy, 303 then compared with human results. Using a common scale, Recall $= \frac{x}{y}$ and Accuracy $= \frac{x}{r}$, where x is the correct answers found, y is the answers given, and r is the total correct answers. Table 1 shows the experimental results. Rows labeled Student/w.BoH, ChatGPT3.5/w.FS, and ChatGPT4/w.FS in the ta-310 ble correspond to the experimental results for the baselines: students with a Bachelor's degree in Hu-311 manities, ChatGPT-3.5 using few-shot examples, 312 and ChatGPT-4 using few-shot examples, respectively. Columns P_1 through P_6 in the table detail 314

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. 315

316

317

318

319

320

321

322

323

324

325

326

328

329

4 InterIDEAS Dataset Overview

In this study, we focus on books originally written 330 or have been translated in English. To date, we 331 have analyzed over 45,000 pages of modern phi-332 losophy available in English. Still expanding, our 333 dataset has amassed over 15,000 cross-referential 334 data pairs, encompassing more than 3,150 philoso-335 phers and philosophical schools, covering the ma-336 jority of both during this period. Our periodization corresponds to the so-called "modern period" in 338 the humanities. Despite its lack of pinpointable timeline, the usual consensus is that the modern 340 period is loosely bound by the beginning of the 341 Industrial Revolution (circa 1760) and the end of 342 WWII (1945). We slightly extended the timeline 343 to address the time lag between historical events and their intellectual stimuli and reactions. In se-345 lecting texts, we balanced coverage with repre-346

	Prompt for Identify Reference	
Static information: You are a professional philosopher. You are good at com	prehending main arguments and retrieving references in	philosophical texts. Let us think step by step.
Questions: Within the passage, please list all the references to exter schools of thoughts. 1. Please limit yourself to explicit external reference 2. Use the author's name/the name of a group to sp "some philosophers claim"), list their authors in order as ' the name of the source. 3. If one external source is mentioned several times	nal textual sources, including specific authors, quotes, b s. ecify each reference and list them separately; for referer Unidentified 1," "Unidentified 2," etc. For collective/unide to enable the current author to make different claims, pl	ooks, ideologies, religions, and literary or philosophical nees whose author is unidentified (like "a poet says," antifiable authorship, such as the Bible, specify them by ease also treat the case as multiple references and list
them separately. 4. If the identified reference includes a reference to by putting an asterisk before it and referring to it as "auth Please do not explain and just give the answer!	another source, please list the second-order reference a or of the first-order reference—author of the second-ord	Iter the first-order one. Signify the second-order reference er reference".
Few-shot instances: One is struck, in the trials of 1782-9, by the increase in te	ension. There is a new severity towards the poor, a conc	erted rejection of evidence, a rise in mutual mistrust,
Answer of Instances: P. Chaunu Input: A few consecutive paragraphs from a book	Refere	nce List
Prompt for Type	Prompt for Intertextual Function	Prompt for Sentiment
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 steo. Questions: For each reference you identified in Previous Question, please describe its content with one of the following descriptions: cominal, verbal, thematic. Few-shot instances: 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)	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. Questions: For each reference you identified in Previous Question, please evaluate the intertextual function it plays Few-shot instances: 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) Answer of Instances:	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. Questions: For each reference you identified in Previous Question, please rate the current work's sentiment toward each reference and characterize the sentiment in terms of negative, neutral, positive Few-shot Instances: , 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)
Answer of Instances: P. Chaunu: Nominal ; Input: A few consecutive paragraphs from a book + Reference List	P. Chaunu: 2. Contextual Explanation; Input: Reference: InterFunction A few consecutive paragraphs from a book + Reference List	Answer of Instances: P. Chaunu: Positive; Input: A few consecutive paragraphs from a book + Reference List

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

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
ChatGPT3.5	0.46	0.58	0.66	0.71	0.67	0.63	0.54	0.61	0.53	0.47	0.25	0.43
ChatGPT3.5/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
ChatGPT4/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

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

4.1 Metadata Format

347

351

355

356

360

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 books, ideologies, religions, historical events, and words and deeds of other people. With our deliberately loose definition guiding LLMs 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.

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

We present a metadata schema specifically designed for the analysis of intertextual references within humanities writing (see Appendix L for detail). The schema facilitates the categorization and detailed examination of references, and their content type, intertextual functions, and sentiment.

376

377

The provided dataset is an organized compilation of bibliographic entries related to philosophy books, encompassing detailed attributes for each book. These attributes include the Book Title, the Reference Name, Linked directly to each reference is the Content Type, which provides detailed information sorted into the nominal, the verbal, and the thematic. This entity captures the essence of each reference through the identification of specified names and titles, presenting quotations from other texts, and giving brief summaries for loose, unspecified discussion of external references, respectively. Each reference is also associated with an Intertextual Function, which describes the role the reference plays in the text-ranging from namedropping (ND) and contextual explanation (CEx) to critical engagement (CEn) and conceptual applica-394 tion or expansion (CAoE). This classification helps us understand the extent of interaction between the current work and the referred content. Furthermore, the Sentiment assesses the current author's sentiment towards each reference, which is categorized as negative, neutral, or positive. This evaluation is 400 crucial for discerning the author's perspective and 401 the reference's intended effect on readers' under-402 standing. The relationships among these values are 403 structured to ensure an one-to-one correspondence 404 between a reference and its content type, intertex-405 tual function, and sentiment. 406

Based on our dataset, Nominal references are 407 the most common, constituting 67.9% of the data, 408 followed by thematic and verbal references. In sen-409 timent analysis, neutral sentiments predominate at 410 77.4%, with positive and negative sentiments at 411 13.1% and 9.5% respectively. For intertextual func-412 tions, name-dropping is most frequent, making up 413 52% of the instances, whereas critical engagement 414 and contextual explanation are also significant, and 415 conceptual application or expansion is relatively 416 rare. These statistics illustrate the dominance of 417 nominal referencing and neutral sentiments in the 418 dataset, with name-dropping being the primary in-419 tertextual function. Meanwhile, authors' attitudes 420 are crucial in determining the depth of their en-421 gagement with others' ideas and actions (see Ap-422 423 pendix J for detailed information). Negative attitudes are often suggested by explicit criticism. 424 People, events, or works that the current authors 425 feel impartial about are usually cursorily discussed. 426 Our general statistics of our dataset also uncover 427

features of modern philosophical writings. First, the dominance of neutral and positive sentiments show that the field is largely organized by amicability. Second, the distribution of sentiments across intertextual function suggests that in constructing philosophical arguments, philosophers generally adapt the style of discussion (recorded as "function" in the dataset) rather than the choice of materials ("type") to reflect their perspectives on individuals, schools of thought, and events ("sentiment"). Comparing the number of positive references with that of the negative ones, we find that philosophers express amicability more overtly and more frequently. They also demonstrate more consensus in their positive acknowledgments of others' work than in their negative critiques. In other words, modern philosophy is primarily organized by amicability rather than confrontation.

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

5 Applications of InterIDEAS in Philosophy and LLMs

5.1 Analysis in InterIDEAS for Philosophy :

Philosophy's canon is vast and densely cross-449 referential. Automating citation extraction lets 450 scholars visualise how ideas propagate across cen-451 turies, schools, and authors-something infeasible 452 by hand at scale. The following section illustrates 453 how our proposed method supports both diachronic 454 and synchronic analyses of philosophical texts in 455 a quantifiable manner. We extract the 50 most fre-456 quent references appeared in at least 3 texts. The 457 word map 3a confirms the interdisciplinary nature 458 of philosophy. Besides acclaimed philosophers 459 and philosophical schools, we find religions (e.g., 460 "Christianity," "God," "Buddha," and "The Bible") 461 and political events and entities (e.g., "Roman Em-462 pire," "British Empire," and "French Revolution") 463 constitutive to philosophical discussion. We ex-464 tract all the individuals from these common refer-465 ences, analysis their life-span and major location 466 of intellectual activity on a timeline and a map re-467 spectively. We finds out that modern philosophers 468 regard ancient, enlightenment, and contemporane-469 ous philosophy in the Mediterranean region and 470 the English Channel region as their shared base 471 of intellectual inquiry. Upon this foundation, we 472 connect writers of antiquity, the Enlightenment era, 473 and the modern era by a mapping network Fig. 3c 474 which tracks the flow of ideas. The network shows, 475 for example, how likely a modern philosopher who 476 has referred to Solon would also be influenced by 477



(a) Word Map

England

478

479

480

481

482

484

485

486

487

488

489

490

491

492

493

494

496

497

498

499

501

504

Boxhorn F

(b) Flow of thought: a graph of referenced philoso(c) Relationship networks for shared references of Emile Faguet and Russell phers from ancient to modern era.

Figure 3: Philosophical references analysis



Figure 4: Pie chart summary for sentiment

Voltaire and moreover by Schopenhauer. Our chart suggests that two important intellectual tradition for modern philosophy are Plato-Rousseau-Hegel and Plato-John Stuart Mill-Engels.

By statistically presenting the proportion of different authors' attitudes in Fig. 4, we identify possible similarities in the tones of their writings. For instance, Georg Jellinek and Franz Oppenheimer may share a more placid style, while Russell's writing tends to be more polemical. Moreover, our dataset allows us to construct intertextual network for any of the two philosophers, as shown by Fig. 3b, to identify previously unknown relationships. In this circumstance, while Bertrand Russell and Émile Faguet are rarely discussed together in philosophical discussion, their shared strong sentiment for Homer and against John Stuart Mill cast 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 Sentiment Classification Enhancement for Language Models

To demonstrate the potential usage of the collected dataset in AI tasks, we create 2,236 referenceattitude pairs from our dataset. 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, respectively.

505

507

508

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

536

537

538

539

540

541

542

543

For validation, 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), AL-BERT (Lan et al., 2019), and BERTweet (Nguyen et al., 2020); 2) RoBERTa (Liu et al., 2019); 3) XLNet (Yang et al., 2019). On the other hand, the four LLMs can be classified into three types too: 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 also study GPT-40 (Achiam et al., 2023), which is the most state-of-the-art (SOTA) LLM, to do direct inference without any extra training. Furthermore, we randomly choose 5 samples of each label from the training set as the few-shot instances for GPT-40. The performance of all PLMs and LLMs pre-trained for text/sequence classification is compared 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 2. For more

Model	Before fine-tuning/few-shot				Afte	er fine-tui	ning/few-	Computational cost			
Widdei	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	-	-	-

Table 2: 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.

clear demonstration, the confusion matrices of each model are shown and analyzed in Appendix K.

544

546

548

549

552

554

555

556

557

558

560

561

562

563

564

565

569

570

573

574

575

577

579

In Table 2, the average performance improvements before and 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. 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. Besides, 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. In conclusion, our corpus positively contributes to helping language models better understand the philosophical context.

Besides, we present confusion matrices of Llama 3 w./w.o. fine-tuning and GPT-4o w./w.o. fewshot learning adopted for sentiment classification in Fig. 5. Before fine-tuning or few-shot learning, all models tend to favor one class and do not consistently choose the Negative class, despite its abundance in the test set. After fine-tuning, models show a marked preference for negative sentiment, indicating improved performance through fine-tuning.



Figure 5: Confusion matrices of Llama 3 w./w.o. finetuning and GPT-40 w./w.o. few-shot learning adopted for sentiment classification.

6 Conclusion

In this article, we introduce InterIDEAS for extracting and evaluating philosophical intertextuality. Enhanced by both LLMs and philosophical expertise, this dataset provides a robust foundation for exploring intellectual structures and dynamics through references. We propose a systematic methodology to categorize, analyze, and interpret complex relationships within and beyond philosophy. InterIDEAS elucidates the intricate ways in which different discourses influence each other, uncovering latent patterns in philosophy that offer insights to both philosophical studies and AI research.

585

586

587

588

589

590

591

592

593

594

595

596

597

598

Limitations

599

619

622

631

634

635

637

639

640

641

642

643

645

647

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 paralleling 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 in the questions, effectively mitigating this issue. 2) 610 Literal-mindedness: The LLM struggles in literary expressions with complex emotions, such as rhetor-611 ical questions and irony. This aspect has seen some 612 improvement through the addition of few-shot in-613 stances. 3) Stereotyping: Faced with specific input 614 information, such as "Hitler," the LLM tends to 615 respond based on its built-in stereotypes with "neg-616 ative" disregarding the author's potentially "neutral" or "positive" stance. 618

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.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, 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 Maffezioli, 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.

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

- 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*.
- 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.

705

706

711

712

714

716

718

721

722

723

724

725

726

727

731

735

738

739

740

741

742

743

744 745

746

747 748

752 753

754

755

- Oliver Li, Mallika Subramanian, Arkadiy Saakyan, Sky CH-Wang, and Smaranda Muresan. 2023. Norm-Dial: 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, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*. 761

762

765

767

768

770

773

774

775

776

779

780

781

782

783

784

- 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

786

787

791

793

796

797

810

811

813

814

815

817

818

821

823

825

827

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, ChatGPT 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. 836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

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: 20-80 883
- Language Proficiency: Participants must be
 884

.1	a	Ψ	т,	,
ti ll ng s. ci if by	fi th g s ns fy id cla 1, ia	ca ne sp Js er ai "bl	at i e ar i e	
ed f 1 g (d on ne a ve	l, j th th on s s l l : r t er t er t er t er t h ti t j n i d u ni	pl e no (thin u recover all a single cover all a si	e foos other a sector l the l the l the property	a o e h a l a e n i s l a u i v n i, p l o
m		tio		(1
i j oi si ati æp	pr n fy ua pt n	it t l ua	IT I h E al g	נו פ וס י

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

Materials Provided

886

887

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

923

924

926

928

929

931

932

933

- 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, 930 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.

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

C.2 **Ouestions**

C.2.1 Q1 for Reference Iden on

eferences to Within the passage, please list a external textual sources, including ific authors. quotes, books, ideologies, religio d literary or philosophical schools of thought the author's name/the name of a group to spec ch reference: for references whose author is u fied (like "a poet says," "some philosophers "), list their authors in order as "Unidentified Unidentified 2," etc. For collective/unident authorship, such as the Bible, specify them name of the source.

C.2.2 Q2 for Content Type

For each reference you identifie se describe its content with one or more of ollowing descriptions: 1. Nominal, meanin e references that explicitly mention names ner authors, books, collections of works, an r schools of thought in the main text; for no references, signal their content by exact nan d in the passage. If there are multiple nomin erences, separate them by colons. E.g., Mar inal (Marx; The Communist Manifesto) 2. V meaning direct quotation of phrases and ser from other sources; for verbal references, s heir content by abbreviated versions of the qu at only keep the first and the last two words of ote, with ellipses in between. If there are m verbal references, separate them by colons. Iarx: verbal ("the history... class struggles") natic, meaning references to others' claims and motifs not through direct quotes but thi paraphrases; for thematic references, please si heir content by a summary in one or two ph hical terms. If there are multiple thematic re es, separate (child labor) them by colons. E.g., Marx: the

C.2.3 Q3 for Intertextual Fu

For each reference identified in pt 1, please evaluate the intertextual function lays by the closet descriptions below. Class e references by "Name-Dropping," "Contex xplanation," "Critical Engagement," or "Conc Application or Expansion." 1. Name-Dro This category is for when the current wor ly mentions

the names of authors, works, or concepts as repre-982 sentative cases of a phenomenon or an argument, without detailed explanations. 2. Contextual Ex-984 planation: Elements of external sources are mentioned and given some exposition to clarify the source's relevance to the author's argument. These 987 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 ar-991 gument, references that intend to exemplify the author's arguments, etc. 3. Critical Engagement: 993 In this category, the current work actively engages 994 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 999 reference. 4. Conceptual Application or Expan-1000 sion: References that fall into this category are not only explained but are also used as a springboard 1002 for further development of the current work.

C.2.4 Q4 for Sentiment

1004

1005

1006

1007

1009

1010

1011

1012

1013

1014

1015

1016

1019

1020

1022

1023

1024

1025 1026

1027

1028

1030

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

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 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). GPT41031with FS and GPT4 with FPEh show consistent high1032performance with averages of 0.702 and 0.815, re-1033spectively, and low variability (SD < 0.08), with1034normal distribution (p-values > 0.650).1035

1036

1054

D.2 Recall

The updated dataset table presents a comprehensive 1037 statistical analysis of performance scores from var-1038 ious groups, including Human Experts, Students 1039 with and without Book of Humanities (BoH), and 1040 different versions of GPT models. The Human Ex-1041 perts group exhibits nearly perfect scores with an 1042 average of 0.988 and a minimal standard deviation 1043 of 0.026, although their scores do not follow a nor-1044 mal distribution. In contrast, the Student groups 1045 show more variability, with averages of 0.718 and 1046 0.552 for Students with BoH and Other Students, 1047 respectively. The GPT models display a progres-1048 sion in performance from GPT3.5 to our approach 1049 with GPT4, where the latter achieves an impressive 1050 average of 0.833 with a standard deviation of 0.053, 1051 showing a more consistent performance (normality p-value = 0.955). 1053

E Interview with Human Experts

We further surveyed human experts about their 1055 opinions on our dataset. All of our human experts, 1056 who are either university professors of philosophy 1057 or PhD students in the humanities, find this dataset 1058 both intriguing and valuable. Representing a bridge 1059 between traditional academic studies and the latest 1060 technological advancements, our application offers 1061 a novel method for integrating these two fields. One 1062 of our interviewees said, "Given the vast scope of 1063 work that no individual could complete in a life-1064 time, the use of language learning models now 1065 makes this formidable task feasible." Another inter-1066 viewee recognized the philosophical implication of 1067 our approach: "Philosophy is a strange field, with 1068 a style of inquiry sometimes behaving like math-1069 ematics and sometimes like literary studies. The seeming incompatibility between the two sets of 1071 assumptions is what keeps me coming back to it, 1072 and this investigation clarifies a lot." One profes-1073 sor was intrigued by how our approach gives con-1074 crete guidance for practical pedagogical tasks like 1075 designing syllabus and creating analytical assign-1076 ments by showing the interrelations among texts. 1077 A PhD student pointed out that the granularity of the information in the dataset is "just right"; the 1079

Table 3: 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

Table 4: 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	

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:



Figure 6: Data format for fine-tuning.

G **Computational Resources**

All data collection processes and fine-tuning experiments are conducted on a server with 8 NVIDIA 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 2, including the model parameter count, the propor-1100 tion of fine-tuned parameters to the total parameter 1101 count, and the time required for 100 epochs of fine-1102 tuning. For details on the fine-tuning parameters, please refer to Table 5.

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

Η **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

¹ BERT:	https:	<pre>//huggingface.co/google-bert/</pre>
bert-base-u	ncased;	
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;
- https://huggingface.co/Mikael110/ Llama 2. 11ama-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.

1092

1093

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1094 1095

1096

1097

Module	Parameter	Parameter description	Value		
	$r_{\rm LoRA}$	The rank of LoRA matrix	8		
	$\alpha_{\rm LoRA}$	Scaling factor of LoRA matrix	32		
LoRA	δ_{LoRA}	Dropout rate	0.1		
			If XLNet: [layer_1, layer_2]		
	$ heta_{ m LoRA}$		elif Llama or Mistral:		
		Modules to be fine tuned	[q_proj, k_proj, v_proj, o_proj,		
		would be the tuned	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		
Fine-tuning	E	Training epoch	100		
	γ	Weight decay	0.01		
	В	Batch size	16		

Table 5: Hyperparameters details.

(70%), validation set (20%), and test set (10%)with the random seed 42 and shuffling. Specially,for BERTweet, the maximal length of each inputsample is truncated to 128 due to the fixed modelinput dimension.

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138 1139

1140

1141

1142

1143

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 5. 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 finetuned. These hyperparameters work together to optimize the application of LoRA in specific models and tasks, balancing computational cost and model performance.

1144

1145

1162

In terms of fine-tuning, the learning rate r is 1146 set to 1e-4, determining the magnitude of updates 1147 to the model parameters at each step. A smaller 1148 learning rate ensures that the model updates its 1149 parameters in small, precise steps, contributing to 1150 a stable and refined training process, reducing the 1151 risk of instability from large parameter changes. 1152 The training epoch E is set to 100 to avoid under-1153 fitting but might lead to over-fitting. To help with 1154 it, the weight decay rate is set to 0.01 by reducing 1155 the size of the model weights at each update. The 1156 batch size B is set to 16 due to both the size of our 1157 proposed sentiment classification dataset and our 1158 hardware limitation. Additionally, the optimizer 1159 is ADAM, and the load accuracy is 32 bit for all 1160 models. 1161

I Prompts for References

We show all the prompts 7, 8, 9, and 10 we used 1163 as following: 1164

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.

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.
 If one external source is mentioned several times to enable the

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 com pilers 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:

 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.

 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.
 Thematic, meaning references to others' claims, ideas, and

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, Lacretelle, Duport, Pastoret, Target, Bergasse: Nominal (Beccaria, Servan, Dupaty, Lacretelle, Duport, Pastoret, Target, Bergasse, Cahiers); Imperialism: Thematic (Alternative to nationalism);

Communism: Thematic (the Cold War);

...

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

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 closet 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.

 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.
 Critical Engagement: In this category, the current work ac-

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: **Context:** The same as the context in Fig. 7. **Answers of instances:** P. Chaunu: 2. Contextual Explanation;

Beccaria, Servan, Dupaty, Lacretelle, Duport, Pastoret, Target, Bergasse: 1.Name-dropping;

Imperialism: 2. Contextual Explanation;

Communism: 1.Name-dropping;

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

Prompt for Reference 4

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:

Please rate the current work's sentiment toward each reference identified in question 1, and characterize the sentiment in terms of negative, neutral, 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." For second-order references, please assess the author's sentiment to the second-order reference, not the sentiment of the first-order reference to the second-order reference. Please base your judgment only on the provided passage.

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.

Few-shot instances:

In these few shot examples, we gave examples for all sentiments. In your application, please select the most appropriate sentiment. You don't have to find traces of all sentiments within a given passage. These are examples for your answer: **Context:**

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

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

J Data Analysis

In this section, we provide additional data analysis in Table 6, Table 8, Table 7 and Figure 11

Table 6: 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	12	119	145
Expansion			

Table 7: Distribution of sentiment types across reference categories.

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

1165



Figure 11: Pie charts showing the distribution of reference types, sentiment types, and intertextual functions.

Table 8: Distribution of sentiment types across intertextual functions.

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

Supplementary Analysis on Sentiment Κ Classification

1168

1169

1195

1196

1197

1198

The confusion matrices of each PLM or LLM is 1170 shown in Figure 12. It can be observed that both 1171 PLMs and LLMs tend to output a specific class, as 1172 seen in the following patterns: Neutral - BERTweet, 1173 RoBERTa, XLNet, Llama 2, GPT-4; Positive -1174 BERT, ALBERT, Llama 3, Mistral, GPT-2. No-1175 tably, none of the models consistently favors the 1176 Negative class, even though Negative samples are 1177 the most abundant in the test set. This tendency 1178 could be attributed to the differences in the pre-1179 training corpora and methods used for each model. 1180 Additionally, LLMs exhibit more moderate biases 1181 compared to PLMs, especially in more recent mod-1182 els like Llama 3, which also has the largest num-1183 ber of parameters. This can be attributed to the 1184 enhanced language understanding capabilities of 1185 LLMs, driven by their larger parameter counts and 1186 more extensive training corpora. Nonetheless, this 1187 highlights a significant issue: even the most ad-1188 vanced language models suffer from severe mode 1189 collapse when directly performing sentiment clas-1190 sification in a philosophical context. Therefore, 1191 the most straightforward approach to enhance a 1192 language model's understanding of philosophical 1193 texts is fine-tuning. 1194

> 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.

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1215

1216

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.

L Metadata format and description

We present the metadata schema specifically designed for the analysis of intertextual references within humanities writing (as mentioned in Sec-1217 tion 4.1) in Fig. 13. 1218



Figure 12: Confusion matrices of each model adopted for sentiment classification before and after fine-tuning or few-shot learning.



Figure 13: Metadata format and description.